BUKTI KORESPONDENSI

- Judul : Implementation of Four-Tier Multiple-Choice Instruments Based on the Partial Credit Model in Evaluating Students' Learning Progress
- Jurnal : European Journal of Educational Research Volume 10, Issue 2, 825 -840. ISSN: 2165-8714 <u>http://www.eu-jer.com/</u> <u>https://doi.org/10.12973/eu-jer.10.2.825</u>



Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Your manuscript ID#2011240749 has been received

1 pesan

European Journal of Educational Research <editor@eu-jer.com> Balas Ke: European Journal of Educational Research <editor@eu-jer.com> Kepada: European Journal of Educational Research <syukrulhamdi@uny.ac.id> 16 Desember 2020 14.10

Dear Dr. Syukrul Hamdi (syukrulhamdi@uny.ac.id),

This mail has been sent automatically by the system.

Your manuscript entitled "Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (ID#2011240749) has been submitted successfully.

The link of your manuscript: https://eu-jer.com/aa/lib/elfinder/files/2011240749/MS_EUJER_ID_2011240749_1.doc

We will inform you about the developments of your paper. Thank you for your interest to our journal.

Best regards.

Editorial Office, European Journal of Educational Research www.eu-jer.com editor@eu-jer.com



Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Corrections request for the manuscript ID# 2011240749

5 pesan

Editor - European Journal of Educational Research <editor@eu-jer.com> Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> 17 Februari 2021 22.01

Dear Dr. Syukrul Hamdi,

After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Implementation of Fourtier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (Manuscript EU-JER ID#2011240749) can be published on condition that corrections are made.

Please consider the reviewers' reports and emendations about your paper, please edit your manuscript and resend the finalized paper via email to us as soon as possible. In addition, we request to fill out the attached correction report what you have done as a word file. Please also highlight the edited parts (or use track changes mode in word).

After we check your manuscript, we will send you the acceptance letter. The deadline for sending your finalized paper is **March 10, 2021** in order to publish in our next issue.

1- A native speaker should check the language of the whole paper as a proofreading lastly.

2- Please check all references for compatibility to APA 7 style (Please see the citation guide page in our web site: https://eu-jer.com/citation-guide).

3- Please try to use at least 2 references from our journal (especially from Vol.9 and Vol.8) in order to increase the impact factor in Scopus.

4- Please provide English translation of the title of non English sources as at the below: Eq.

Bussieres, E.-L., St-Germain, A., Dube, M., & Richard, M.-C. (2017). Efficacite et efficience des programmes de transition a la vie adulte: Une revue systematique [Effectiveness and efficiency of adult transition programs: A systematic review]. *Canadian Psychology/ Psychologie Canadienne, 58*(1), 354–365. https://doi.org/10.1037/cap0000104

Note: If all of the corrections don't be completed, the paper will be refused. If you object to any correction, please explain this in your correction report.

Please confirm when you get this email. We are looking forward to hearing you.

Best regards,

Ahmet Savas, Ph.D. Editor-in-Chief, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

8 lampiran

- EU-JER_REVIEWER_FORM_R2614.docx 336K
- EU-JER_REVIEWER_FORM_R2615.docx
- MS_EUJER_ID_2011240749_R2611.doc 918K
- MS_EUJER_ID_2011240749_R2612.doc 928K
- MS_EUJER_ID_2011240749_R2614.doc 918K

W	CORRECTION RE	PORT_EU-JER.	docx
	19K		
	19K		

EU-JER_REVIEWER_FORM_R2611.docx 135K

W	EU-JER	_REVIEWER_	FORM	_R2612.docx
E	136K			

Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> Kepada: Lukman Abdul Rauf Laliyo <lukman.laliyo@ung.ac.id> 18 Februari 2021 02.29

------ Forwarded message ------Dari: Editor - European Journal of Educational Research <editor@eu-jer.com> Date: Rab, 17 Feb 2021 pukul 22.02 Subject: Corrections request for the manuscript ID# 2011240749 To: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Dear Dr. Syukrul Hamdi,

After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Implementation of Fourtier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (Manuscript EU-JER ID#2011240749) can be published on condition that corrections are made.

Please consider the reviewers' reports and emendations about your paper, please edit your manuscript and resend the finalized paper via email to us as soon as possible. In addition, we request to fill out the attached correction report what you have done as a word file. Please also highlight the edited parts (or use track changes mode in word).

After we check your manuscript, we will send you the acceptance letter. The deadline for sending your finalized paper is **March 10, 2021** in order to publish in our next issue.

1- A native speaker should check the language of the whole paper as a proofreading lastly.

2- Please check all references for compatibility to APA 7 style (Please see the citation guide page in our web site: https://eu-jer.com/citation-guide).

3- Please try to use at least 2 references from our journal (especially from Vol.9 and Vol.8) in order to increase the impact factor in Scopus.

4- Please provide English translation of the title of non English sources as at the below:

Eg.

Bussieres, E.-L., St-Germain, A., Dube, M., & Richard, M.-C. (2017). Efficacite et efficience des programmes de transition a la vie adulte: Une revue systematique [Effectiveness and efficiency of adult transition programs: A systematic review]. *Canadian Psychology/ Psychologie Canadienne, 58*(1), 354–365. https://doi.org/10.1037/cap0000104

Note: If all of the corrections don't be completed, the paper will be refused. If you object to any correction, please explain this in your correction report.

Please confirm when you get this email. We are looking forward to hearing you.

Best regards,

Ahmet Savas, Ph.D. Editor-in-Chief, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

8 lampiran

EU-JER_REVIEWER_FORM_R2614.docx 336K

EU-JER_REVIEWER_FORM_R2615.docx

MS_EUJER_ID_2011240749_R2612.doc 928K
MS_EUJER_ID_2011240749_R2614.doc 918K
CORRECTION REPORT_EU-JER.docx
EU-JER_REVIEWER_FORM_R2611.docx 135K
EU-JER_REVIEWER_FORM_R2612.docx 136K

Editor - European Journal of Educational Research <editor@eu-jer.com> Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> 20 Februari 2021 19.19

Dear Dr. Hamdi,

As you may remember, we asked you to cite our journal in your article before in order to increase our impact factor. However, we have now seen that we have achieved the impact value we wanted in Scopus and we gladly saw that our citation count exceeded the desired limit value. As the journal management, we have decided not to cite our journal in the articles published in our journal since the next issues. In this case, we request that you do not cite our journal in your article and delete the citations and references related our journal, if any. I am sorry about the inconsistency.

We are looking forward to getting your revised paper.

Best regards,

Ahmet C. Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

On 2/17/2021 6:01 PM, Editor - European Journal of Educational Research wrote:

Dear Dr. Syukrul Hamdi,

After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (Manuscript EU-JER ID#2011240749) can be published on condition that corrections are made.

Please consider the reviewers' reports and emendations about your paper, please edit your manuscript and resend the finalized paper via email to us as soon as possible. In addition, we request to fill out the attached correction report what you have done as a word file. Please also highlight the edited parts (or use track changes mode in word).

After we check your manuscript, we will send you the acceptance letter. The deadline for sending your finalized paper is **March 10, 2021** in order to publish in our next issue.

1- A native speaker should check the language of the whole paper as a proofreading lastly.

2- Please check all references for compatibility to APA 7 style (Please see the citation guide page in our web site: https://eu-jer.com/citation-guide).

3- Please try to use at least 2 references from our journal (especially from Vol.9 and Vol.8) in order to increase the impact factor in Scopus.

4- Please provide English translation of the title of non English sources as at the below: Eg.

Bussieres, E.-L., St-Germain, A., Dube, M., & Richard, M.-C. (2017). Efficacite et efficience des programmes de transition a la vie adulte: Une revue systematique [Effectiveness and efficiency of adult transition programs: A systematic review]. *Canadian Psychology/ Psychologie Canadienne, 58*(1), 354–365. https://doi.org/10.1037/cap0000104 Note: If all of the corrections don't be completed, the paper will be refused. If you object to any correction, please explain this in your correction report.

Please confirm when you get this email. We are looking forward to hearing you.

Best regards,

Ahmet Savas, Ph.D. Editor-in-Chief, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> Kepada: Editor - European Journal of Educational Research <editor@eu-jer.com>

20 Februari 2021 22.26

Dear. Editor of EU-JER

Thank you for the information. We will immediately revise the article according to the suggestions of the reviewers. Thank you

Best regards Syukrul Hamdi

Pada tanggal Sab, 20 Feb 2021 pukul 19.20 Editor - European Journal of Educational Research <editor@eu-jer.com> menulis:

Dear Dr. Hamdi,

As you may remember, we asked you to cite our journal in your article before in order to increase our impact factor. However, we have now seen that we have achieved the impact value we wanted in Scopus and we gladly saw that our citation count exceeded the desired limit value. As the journal management, we have decided not to cite our journal in the articles published in our journal since the next issues. In this case, we request that you do not cite our journal in your article and delete the citations and references related our journal, if any. I am sorry about the inconsistency.

We are looking forward to getting your revised paper.

Best regards,

Ahmet C. Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

On 2/17/2021 6:01 PM, Editor - European Journal of Educational Research wrote:

Dear Dr. Syukrul Hamdi,

After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (Manuscript EU-JER ID#2011240749) can be published on condition that corrections are made.

Please consider the reviewers' reports and emendations about your paper, please edit your manuscript and resend the finalized paper via email to us as soon as possible. In addition, we request to fill out the attached correction report what you have done as a word file. Please also highlight the edited parts (or use track changes mode in word).

After we check your manuscript, we will send you the acceptance letter. The deadline for sending your finalized paper is **March 10, 2021** in order to publish in our next issue.

1- A native speaker should check the language of the whole paper as a proofreading lastly.

2- Please check all references for compatibility to APA 7 style (Please see the citation guide page in our

web site: https://eu-jer.com/citation-guide).
3- Please try to use at least 2 references from our journal (especially from Vol.9 and Vol.8) in order to increase the impact factor in Scopus.
4- Please provide English translation of the title of non English sources as at the below:
Eg.
Bussieres, E.-L., St-Germain, A., Dube, M., & Richard, M.-C. (2017). Efficacite et efficience des programmes de transition a la vie adulte: Une revue systematique [Effectiveness and efficiency of adult transition programs: A systematic review]. *Canadian Psychology/ Psychologie Canadienne, 58*(1), 354–365. https://doi.org/10.1037/cap0000104

Note: If all of the corrections don't be completed, the paper will be refused. If you object to any correction, please explain this in your correction report.

Please confirm when you get this email. We are looking forward to hearing you.

Best regards,

Ahmet Savas, Ph.D. Editor-in-Chief, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

Editor - European Journal of Educational Research <editor@eu-jer.com> Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Dear Dr. Hamdi ,

Thank you for your kind reply.

We are looking forward to getting your revised paper.

Best regards,

Ahmet C. Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

On 2/20/2021 6:26 PM, Syukrul Hamdi UNY wrote:

Dear. Editor of EU-JER

Thank you for the information. We will immediately revise the article according to the suggestions of the reviewers. Thank you

Best regards Syukrul Hamdi

Pada tanggal Sab, 20 Feb 2021 pukul 19.20 Editor - European Journal of Educational Research <editor@eu-jer.com> menulis:

Dear Dr. Hamdi,

As you may remember, we asked you to cite our journal in your article before in order to increase our impact factor. However, we have now seen that we have achieved the impact value we wanted in Scopus and we gladly saw that our citation count exceeded the desired limit value. As the journal management, we have decided not to cite our journal in the articles published in our journal since the next issues. In this case, we request that you do not cite our journal in your article and delete the citations and references related our journal, if any. I am sorry about the inconsistency.

We are looking forward to getting your revised paper.

Best regards,

21 Februari 2021 00.12

Ahmet C. Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com
On 2/17/2021 6:01 PM, Editor - European Journal of Educational Research wrote:
Dear Dr. Syukrul Hamdi,
After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (Manuscript EU-JER ID#2011240749) can be published on condition that corrections are made.
Please consider the reviewers' reports and emendations about your paper, please edit your manuscript and resend the finalized paper via email to us as soon as possible. In addition, we request to fill out the attached correction report what you have done as a word file. Please also highlight the edited parts (or use track changes mode in word).
After we check your manuscript, we will send you the acceptance letter. The deadline for sending your finalized paper is March 10, 2021 in order to publish in our next issue.
 1- A native speaker should check the language of the whole paper as a proofreading lastly. 2- Please check all references for compatibility to APA 7 style (Please see the citation guide page in our web site: https://eu-jer.com/citation-guide). 3- Please try to use at least 2 references from our journal (especially from Vol.9 and Vol.8) in order to increase the impact factor in Scopus. 4- Please provide English translation of the title of non English sources as at the below: Eq.
Bussieres, EL., St-Germain, A., Dube, M., & Richard, MC. (2017). Efficacite et efficience des programmes de transition a la vie adulte: Une revue systematique [Effectiveness and efficiency of adult transition programs: A systematic review]. <i>Canadian Psychology/ Psychologie Canadienne, 58</i> (1), 354–365. https://doi.org/10.1037/cap0000104
Note: If all of the corrections don't be completed, the paper will be refused. If you object to any correction, please explain this in your correction report.
Please confirm when you get this email. We are looking forward to hearing you.
Best regards,
Ahmet Savas, Ph.D. Editor-in-Chief, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya. (To support the "Green UNY movement", it is recommended not to print the contents of this email and its

(To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments)

Universitas Negeri Yogyakarta

www.uny.ac.id



Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

GENTLE REMINDER: Corrections request for the manuscript ID# 2011240749

4 pesan

Editor - European Journal of Educational Research <editor.eujer@gmail.com> Balas Ke: editor@eu-jer.com Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> 10 Maret 2021 15.10

GENTLE REMINDER

----- Forwarded Message ------

Subject:Corrections request for the manuscript ID# 2011240749 Date:Wed, 17 Feb 2021 18:01:43 +0300 From:Editor - European Journal of Educational Research <editor@eu-jer.com> To:Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Dear Dr. Syukrul Hamdi,

After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Implementation of Fourtier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (Manuscript EU-JER ID#2011240749) can be published on condition that corrections are made.

Please consider the reviewers' reports and emendations about your paper, please edit your manuscript and resend the finalized paper via email to us as soon as possible. In addition, we request to fill out the attached correction report what you have done as a word file. Please also highlight the edited parts (or use track changes mode in word).

After we check your manuscript, we will send you the acceptance letter. The deadline for sending your finalized paper is **March 10, 2021** in order to publish in our next issue.

1- A native speaker should check the language of the whole paper as a proofreading lastly.

2- Please check all references for compatibility to APA 7 style (Please see the citation guide page in our web site: https://eu-jer.com/citation-guide).

3- Please try to use at least 2 references from our journal (especially from Vol.9 and Vol.8) in order to increase the impact factor in Scopus.

4- Please provide English translation of the title of non English sources as at the below: Eg.

Bussieres, E.-L., St-Germain, A., Dube, M., & Richard, M.-C. (2017). Efficacite et efficience des programmes de transition a la vie adulte: Une revue systematique [Effectiveness and efficiency of adult transition programs: A systematic review]. *Canadian Psychology/ Psychologie Canadienne, 58*(1), 354–365. https://doi.org/10.1037/cap0000104

Note: If all of the corrections don't be completed, the paper will be refused. If you object to any correction, please explain this in your correction report.

Please confirm when you get this email. We are looking forward to hearing you.

Best regards,

Ahmet Savas, Ph.D. Editor-in-Chief, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

8 lampiran

EU-JER_REVIEWER_FORM_R2614.docx

336K

EU-JER_REVIEWER_FORM_R2615.docx 339K MS_EUJER_ID_2011240749_R2611.doc W 918K MS_EUJER_ID_2011240749_R2612.doc W 928K MS_EUJER_ID_2011240749_R2614.doc W 918K CORRECTION REPORT_EU-JER.docx W 19K EU-JER_REVIEWER_FORM_R2611.docx W 135K EU-JER_REVIEWER_FORM_R2612.docx 136K

Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> Kepada: European Journal of Educational Research <editor@eu-jer.com> 10 Maret 2021 19.57

Dear, Editor of EUJER

I have recently revised the article according to the suggestions and input from reviewers R2611, R2612, R2614 and R2615. I attached the result of the revised article and correction report. Thank you

Best regards Syukrul Hamdi [Kutipan teks disembunyikan]

2 lampiran

CORRECTION REPORT_EU-JER.docx

MS_EUJER_ID_2011240749_Revised.doc 1897K

Editor - European Journal of Educational Research <editor@eu-jer.com> Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> 10 Maret 2021 21.34

Dear Dr. Hamdi,

We have received your revised paper and correction report. We have sent them to our reviewers again in order to check. We will inform you when we get the result from our reviewers.

If the reviewers confirm your revised paper, we will send the acceptance letter to you.

Thank you for your patience.

Best regards,

Ahmet Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com [Kutipan teks disembunyikan]

[Kutipan teks disembunyikan]

Email Universitas Negeri Yogyakarta - GENTLE REMINDER: Corrections request for the manuscript ID# 2011240749

Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya. (To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments)

Universitas Negeri Yogyakarta www.uny.ac.id

Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> Kepada: Editor - European Journal of Educational Research <editor@eu-jer.com> 11 Maret 2021 04.01

Dear, Editor of EUJER

Thank you very much

Best regards Syukrul Hamdi

[Kutipan teks disembunyikan]



Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Second round corrections request for the manuscript ID# 2011240749

1 pesan

Editor - European Journal of Educational Research <editor@eu-jer.com> Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> 18 Maret 2021 18.13

Dear Dr. Hamdi,

Please see the attached file as the second round corrections.

Please remove the old highlights and re-highlight for new edited parts. We don't need a new correction report.

We are looking forward to getting your second revised paper until March 25, 2021.

Best regards,

Ahmet Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com

On 11-Mar-21 12:01 AM, Syukrul Hamdi UNY wrote:

Dear, Editor of EUJER

Thank you very much

Best regards Syukrul Hamdi

Pada tanggal Rab, 10 Mar 2021 pukul 21.34 Editor - European Journal of Educational Research <<u>editor@eu-jer.com</u>> menulis:

Dear Dr. Hamdi,

We have received your revised paper and correction report. We have sent them to our reviewers again in order to check. We will inform you when we get the result from our reviewers.

If the reviewers confirm your revised paper, we will send the acceptance letter to you.

Thank you for your patience.

Best regards,

Ahmet Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

On 10-Mar-21 3:57 PM, Syukrul Hamdi UNY wrote:

Dear, Editor of EUJER

I have recently revised the article according to the suggestions and input from reviewers R2611, R2612, R2614 and R2615. I attached the result of the revised article and correction report. Thank you

Best regards Syukrul Hamdi Pada tanggal Rab, 10 Mar 2021 pukul 15.10 Editor - European Journal of Educational Research <editor.eujer@gmail.com> menulis:

GENTLE REMINDER

----- Forwarded Message ------

Subject:Corrections request for the manuscript ID# 2011240749 Date:Wed, 17 Feb 2021 18:01:43 +0300

From:Editor - European Journal of Educational Research <editor@eu-jer.com> To:Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Dear Dr. Syukrul Hamdi,

After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (Manuscript EU-JER ID#2011240749) can be published on condition that corrections are made.

Please consider the reviewers' reports and emendations about your paper, please edit your manuscript and resend the finalized paper via email to us as soon as possible. In addition, we request to fill out the attached correction report what you have done as a word file. Please also highlight the edited parts (or use track changes mode in word).

After we check your manuscript, we will send you the acceptance letter. The deadline for sending your finalized paper is **March 10, 2021** in order to publish in our next issue.

1- A native speaker should check the language of the whole paper as a proofreading lastly.

2- Please check all references for compatibility to APA 7 style (Please see the citation guide page in our web site: https://eu-jer.com/citation-guide).

3- Please try to use at least 2 references from our journal (especially from Vol.9 and Vol.8) in order to increase the impact factor in Scopus.

4- Please provide English translation of the title of non English sources as at the below: Eg.

Bussieres, E.-L., St-Germain, A., Dube, M., & Richard, M.-C. (2017). Efficacite et efficience des programmes de transition a la vie adulte: Une revue systematique [Effectiveness and efficiency of adult transition programs: A systematic review]. *Canadian Psychology/ Psychologie Canadienne, 58*(1), 354–365. https://doi.org/10.1037/cap0000104

Note: If all of the corrections don't be completed, the paper will be refused. If you object to any correction, please explain this in your correction report.

Please confirm when you get this email. We are looking forward to hearing you.

Best regards,

Ahmet Savas, Ph.D. Editor-in-Chief, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya.

(To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments)

Email Universitas Negeri Yogyakarta - Second round corrections request for the manuscript ID# 2011240749

Universitas Negeri Yogyakarta www.uny.ac.id

Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya. (To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments)

Universitas Negeri Yogyakarta www.uny.ac.id

www.urry.ao.id

mage: 2nd ROUND_MS_EUJER_ID_2011240749_.doc 1913K



Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Acceptance Letter for the Manuscript ID#2011240749

2 pesan

Editor - European Journal of Educational Research <editor@eu-jer.com> Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> 22 Maret 2021 14.02

Dear Dr. Syukrul Hamdi,

Congratulation! After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (ID#2011240749) has been accepted. It is scheduled for publication in the Volume 10 Issue 2 of the "European Journal of Educational Research".

We kindly ask you to pay the article processing fee USD 500 and USD 100 transaction fee + tax of our bank (totally USD 600) via bank wire transfer. Kindly acknowledge invoice of this acceptance letter. Payment due date: March 26, 2021.

BANK WIRE TRANSFER INFORMATION : NAME OF BENEFICIARY: Ahmet Cezmi SAVAŞ ADDRESS OF BENEFICIARY: Degirmicem District Ozgurluk Str. No:32B , Zipcode:27090, Gaziantep, TURKEY PHONE OF BENEFICIARY: +90 (342) 909 61 90 CORRESPONDENT BANK CHARGER: REMITTER AMOUNT: USD 600 PAYMENT DETAIL: EU-JER_Manuscript ID#2011240749 BANK NAME: QNB Finansbank BANK ADDRESS: Esentepe Mahallesi Büyükdere Caddesi Kristal Kule Binası No:215 Şişli - İstanbul BRANCH OF THE BANK: ENPARA BRANCH CODE: 3663 ACCOUNT NUMBER: 88177946 IBAN: TR66 0011 1000 0000 0088 1779 46 SWIFT CODE: FNNBTRISXXX

Or you can use wise.com for the money transfer to our Turkish Lira account easily. The amount will be 4680 TL, the IBAN of TL is TR45 0011 1000 0000 0082 7703 80, the name of account holder is "Ahmet Cezmi SAVAŞ".

After payment, we will send the gallery proof of your paper. The galley proofs must be returned to us within 2 calendar days. Furthermore, you are responsible for any error in the published paper due to your oversight.

Thank you very much for submitting your article to the journal of "European Journal of Educational Research". We welcome your contributions in future.

Best regards.

Ahmet C. Savas Ph.D. Editor, European Journal of Educational Research http://www.eu-jer.com editor@eu-jer.com

On 20-Mar-21 6:44 PM, Syukrul Hamdi UNY wrote:

Dear, Editor of EUJER

I have recently revised the article according to the suggestions from editor. I attached the result of the revised article. Thank you

Best regards Syukrul Hamdi

Pada tanggal Kam, 18 Mar 2021 pukul 18.13 Editor - European Journal of Educational Research <editor@eu-jer.com> menulis:

Dear Dr. Hamdi,

Please see the attached file as the second round corrections.

Please remove the old highlights and re-highlight for new edited parts. We don't need a new correction report.

We are looking forward to getting your second revised paper until March 25, 2021.

Best regards,

Ahmet Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com

On 11-Mar-21 12:01 AM, Syukrul Hamdi UNY wrote:

Dear, Editor of EUJER

Thank you very much

Best regards Syukrul Hamdi

Pada tanggal Rab, 10 Mar 2021 pukul 21.34 Editor - European Journal of Educational Research <editor@eu-jer.com> menulis:

Dear Dr. Hamdi,

We have received your revised paper and correction report. We have sent them to our reviewers again in order to check. We will inform you when we get the result from our reviewers.

If the reviewers confirm your revised paper, we will send the acceptance letter to you.

Thank you for your patience.

Best regards,

Ahmet Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

On 10-Mar-21 3:57 PM, Syukrul Hamdi UNY wrote:

Dear, Editor of EUJER

I have recently revised the article according to the suggestions and input from reviewers R2611, R2612, R2614 and R2615. I attached the result of the revised article and correction report. Thank you

Best regards Syukrul Hamdi

Pada tanggal Rab, 10 Mar 2021 pukul 15.10 Editor - European Journal of Educational Research <editor.eujer@gmail.com> menulis:

GENTLE REMINDER



Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya.

10/10/21 22.24	15/	10/	21	22.	24
----------------	-----	-----	----	-----	----

Email Universitas Negeri	Yoqvakarta - Acce	eptance Letter for t	he Manuscrip	t ID#2011240749
Email emilie english				

(To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments) Universitas Negeri Yogyakarta

www.uny.ac.id

Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya. (To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments) Universitas Negeri Yogyakarta

www.uny.ac.id

Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya. (To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments) Universitas Negeri Yogyakarta

www.uny.ac.id

Acceptance Letter for the EU-JER_Manuscript_ID#2011240749.pdf 863K

Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> Kepada: Editor - European Journal of Educational Research <editor@eu-jer.com> 23 Maret 2021 11.58

Dear, Editor of EU-JER

We are glad to receive the information regarding our article which is accepted for further publication. I inform you that I have already processed a payment through Bank Negara Indonesia for around USD 600. The staff of bank said that it will arrive at the Destination Bank in about 1-2 days. Thank you

Best regards,

Syukrul Hamdi [Kutipan teks disembunyikan]

Proof of Payment for publication of EU-JER_USD 600.pdf

Bukti Submit di EuJER

$\leftrightarrow \ \ \rightarrow \ \ G$	eu-jer.com/pr	ofile							o , ⊑	☆	*
		Europea	1 Journ	ial o	f Educ	ational R	lesearch				
	Manusc	ript Submissior	n System							٦	
	🕒 start a N	ew submission	SUBMISSIONS	≜ S	YUKRUL HAMDI	Ů SIGN OUT					
	ld	Title	File Cr	reated	Abstract	Author(s)	Review Files	Status			
	2011240749	Implementation of Four- tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress	% 20 16 07	020-12- 5 7:10:21	CLICK TO READ	DR. SYUKRU DR. LUKMAN J MR. MASRI DR. MARIO A DR. CITRA PA	L HAMDI A.R. JALIYO D FIKOLI IBDUILAH NIGORO	Published	C UPDATE		
← → C	a mail.google.co	m/mail/u/0/#search/submit+edit	or%40eu-jer.com/	'FMfcgxwKjv	wxRGZXSdSQPdFo	:DGBtBdMhM			Q	☆	*
= M	Gmail	Q submit editor@eu-jer.co	om			× 幸		Aktif 🔻 🕐	۰.		1 /
r Email		← ▣ ① ₪	⊴ © ⊘.		:				14 dari 34	<	>
Kotak Ma	asuk 664	Your manuscrip	t ID#2011240)749 has	been receiv	ed Kotak Masuk x					Z
 Ditunda Ditunda Terkirim 	ing	European Journal of Edu	icational Research <	<editor@eu-jer< td=""><td>.com></td><td></td><td></td><td>Rab, 16 Des 202</td><td>0 14.10 🟠</td><td>÷</td><td>:</td></editor@eu-jer<>	.com>			Rab, 16 Des 202	0 14.10 🟠	÷	:
D Drof	31	☆A Inggris - > In	donesia 👻 Terjema	ahkan pesan				No	naktifkan untuk:	Inggris	×
r Chat	+	Dear Dr. Syukrul Hamdi (s)	ukrulhamdi@uny.ac.ie	<u>d),</u>							
F		This mail has been sent au	tomatically by the syst	tem.							
Tidak ada i	percakapan	Your manuscript entitled "In	nplementation of Four	-tier Instrume	nts Based on the Ra	sch Model in Evaluating Stud	ents' Learning Progress" (ID#2)	011240749) has bee	n <mark>submitted</mark> suc	cessfull	у.
Mulai	i chat	We will inform you about th	e developments of yo	ur paper. Tha	nk you for your intere	est to our journal.	<u></u>				
- Ruang	+	Best regards.									
Belum a Buat atau ter	ida ruang mukan ruang	Editorial Office, European www.eu-jer.com editor@eu-jer.com	Journal of Educationa	Il Research							

= 附 Gmail		Q eujer	× 3‡	● Aktif ▼ ⑦ 🔅 🏭		
r Email	1			5 dari 13	< >	81
🔲 Kotak Masuk	664					
🕁 Berbintang		Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id></syukrulhamdi@uny.ac.id>		📼 Jum, 16 Okt 2020 03.07 🛛 🛧 🚽	← :	
() Ditunda		kepada European -				~
▶ Terkirim		Dear, Editor of EUJER				0
🗅 Draf	30	I have recently revised the article according to the suggestions. The	ank you			
r Chat	+	Best regards Syukrul Hamdi				_
						+
Tidak ada percakapan Mulai chat		European Journal of Educational Research Research				
r Ruang	+	Bellevine and the second secon				



Implementation of Four-tier *multiple-choice* Instruments Based on the Rasch Model in Evaluating Students' Learning Progress

Abstract: The purpose of this study was to developing and implementation four-tier multiple-choice (hereinafter, 4TMC) instrument with Rasch model to evaluate students' learning progress in explaining the concept of change of state of matter. The data were obtained through development and validation techniques on 20 4TMC items (distributed) to 427 students. On each item, the study applied diagnostic-summative assessment and certainty response index. The students' conceptual understanding level was categorized based on the combination their answer choices; the measurement generated Partial-Credit polytomous Rasch model (data). The data were further processed by WINSTEPS version 4.5.3 software to equate the data interval rate. Analysis of differences based on class level of students using Analysis of Variants (One-way ANOVA). The result revealed that the integration of 4TMC test and Rasch modeling was effective to be treated as the instrument to measure students' learning progress. One-way ANOVA test indicates a difference among the students' competence based on the academic level, in which F_{count} (6, 0142442) > F_{table} (2,39308); df = 422; p <0.05. On top of that, it was discovered that low-ability students see very slow progress due to the lack of knowledge as well as a misconception in explaining the concept as mentioned above.

Commented [L1]: tested

Commented [L2]: what are the characteristics of the item being measured?

Keywords: Learning progress. four-tier, change of state of matter, Rasch model.

Introduction

Central to the notion of science learning is the development of students' scientific understanding of basic concepts of sciences (Hadenfeldt et al., 2013), particularly, change of state of matter (Emden et al., 2018). Aside from the issue, several studies have also highlighted the students' inability to provide an epistemological explanation of basic concepts of sciences (Chi et al., 2018). Efforts to solve the issues, however, have shown little progress, as the students might have more complex perceptions regarding the alternative concept they understand (Morell et al., 2017).

Education practitioners have recommended the utilization of learning progress concept as the instructional method to provide guidance and direction and to adjust the curriculum with the

Commented [L3]: the research results are adjusted to the research objectives (1) development results (2) evaluation results learning process and assessment (Claesgens *et al.*, 2009; Duncan & Hmelo-Silver, 2009; Rogat *et al.*, 2011). Learning progress is defined as a sophisticated and systematic way of thinking. This method is applicable for a learning process, in which the students will undergo gradual progress when learning a topic in a long duration. Its effectiveness is highly dependent on the learning process and the students' learning experience (Duschl et al., 2011). The concept involves certain sets of gradual levels that represent conceptual understanding, ranging from low level up to comprehensive level.

The notion of learning progress is highly distinctive to each student and is dependent to one's learning experience (Rogat et al., 2011); therefore, there is no learning roadmap that is suitable for all kinds of students (Smith et al., 2006). Each student constructs one's understanding in a different way; moreover, the construction process is varied depending on the students' conceptual understanding level (Aktan, 2013). This is to say that each student undergoes a different rate of learning progress, understanding level, and knowledge construction. Simply put, the development of scientific comprehension among students is not linear (Neumann et al., 2013). Therefore, this study regards each level of students' conceptual understanding for more advanced level of understanding (Hadenfeldt et al., 2013). A student who faces difficulty in a certain level of understanding will see a lack of progress to a more advanced level. This in turn hinders the student's ability to construct an epistemological explanation on the basic concepts of science.

Efforts to diagnose the epistemological problems, as mentioned previously, are feasible to conduct if the extent of students' conceptual understanding is formulated. Within this context, the learning progress is treated as the method to evaluate students' conceptual understanding. The diagnostic information generated is reliable to be treated as a reference for the teachers in developing accurate and valid instructional components to guide the students to progress to the next level. Despite the potentials, this study deems that it is challenging for the teachers to construct such an accurate instrument.

Among the diagnostic instruments that are considered applicable is the 4TMC instrument. It is the development of two-tier multiple-choice test recommended by Treagust (1988) and Chandrasegaran et al., (2007). The use of two-tier instrument is familiar in identifying students' understanding in select topics such as electrochemistry (Lu & Bi, 2016), covalent bond (Peterson, Treagust, & Garnett, 1989), and chemical equilibrium (Tyson et al., 1999). Despite its reputation in academia, the two-tier test has raised criticism due to its sole focus on the facts and negligence towards students' understanding (Klassen, 2006). Therefore, several experts propose the renewed version of the test by adding distractor answer choices to strengthen the diagnostic value of the items (Herrmann-Abell & DeBoer, 2011; Sadler, 1998). In addition, some have highlighted the test's weakness in cases where students' tended to pick the answer choice and the reasoning randomly. This illustrates that the students were uncertain and possessed several misconceptions in the first tier question. In such cases, teachers faced difficulty in differentiating between guessed answers and misconceptions (Habiddin & Page, 2019; Hasan et al., 1999).

The criticism laid against the model has sparked the innovation of three-tier and four-tiers instruments. Both instruments feature two multi-level questions, also similar with two-tier test. In the three-tier test, however, the measurement of students' certainty level is conducted simultaneously in both first and second-tier questions; in the meantime, the measurement is conducted separately in the first two tiers (Caleon & Subramaniam, 2010). The value of students' certainty rate ranges from one (very uncertain) to five (very certain).

Three-tier test lacks validity in measuring the students' certainty rate regarding both the answer choice and the reasoning, whether or not the value of certainty rate refers only to the answer choice, to the reasoning, or both. Such weakness will in turn obstructs the evaluation and classification process of students' responses (Arslan et al., 2012). In the four-tier instrument, the measurement of certainty rate also involves the answer choice in the first tier and the reasoning in the third tier (Arslan et al., 2012; Loh et al., 2014). Regarding this

Commented [L4]: need explanation first and then acronym (four-tier multiple-choice) feature, four-tier test is considered more accurate than the three-tier test. Students who pick wrong answer choices with high certainty indicate that they have a very high misconception on the measured item (Hoe & Subramaniam, 2016).

Four-tier instruments are used in studies discussing topics such as physics education (Caleon & Subramaniam, 2010), chemical thermodynamics (Sreenivasulu & Subramaniam, 2013), transition metal (Sreenivasulu & Subramaniam, 2014), acid-base reaction (Hoe & Subramaniam, 2016), and chemical kinetics (Habiddin & Page, 2019). However, it is worth noticing that studies on chemistry topic which employ four-tiers instruments tend to focus on describing alternative conception. To put it another way, the higher the certainty rate is, the stronger the students' alternative conception will be. Despite its potentials, the scholarly discussion has overlooked the implementation of a four-tier diagnostic instrument to measure students' learning progress. Therefore, further analysis is essential on the application of 4TMC test in several domains analyzes by Rasch model approach.

The use of Rasch model has been introduced since the 2000s in the science education research; it features the instrument that integrates diagnostic assessment and summative assessment (Liu, 2012; Wei et al., 2012). On top of that, the diagnostic assessment approach is introduced to conduct an in-depth analysis of the construction process of students' conceptual understanding (Claesgens et al., 2009; Hadenfeldt et al., 2013; Lu & Bi, 2016). This study employs 4MTC test and Rasch modeling as a diagnostic tool to evaluate students' learning progress in explaining the change of state of matter. The study focus revolves around two research questions: 1) How is the effectiveness of 4TMC instrument to evaluate the students' learning progress in explaining concepts of change of state of matter. 2) How is the learning progress in students ranging from the senior high school level up to the senior (fourth) year of college in explaining the concepts?

Commented [L5]: synchronize with the research objectives on the abstract

Methodology

Development Model

This research used a development research referring to the test development model from Wilson. Wilson (2005, 2008) introduces four steps of measurement instrument development: The first step is to the learning progress variable focused on a characteristic measured at a particular time unit. The second step comprises the design process of items or tasks used to measure students' responses. Moreover, the third step involves outcome space, in which the students' responses are categorized into all items related with the learning progress variable. On top of that, the fourth step employs measurement model, such as Rasch model. This recommendation is proven valid to be implemented in developing measurement instrument for different construct variables (Barbera, 2013; Chi et al., 2018; Hadenfeldt et al., 2013; Laliyo, Botutihe, & Panigoro, 2019; Lu & Bi, 2016; Wei et al., 2012; Wilson, 2009; Wind, Tsai, Grajeda, & Bergin, 2018). The study conducted development of measurement instruments by referring to Wilson's recommendation (2005, 2008) and adopted Treagust's framework (1988) of item development. The present study also included two questions related to certainty rate (Arslan et al., 2012; Habiddin & Page, 2019: Hasan et al., 1999). The obtained data were analyzed by Rasch model approach.

Construct Map: Determining Level of Understanding

The first step was to develop the construct of measured variables. The study involved four concepts of change of state of matter: liquid-gas (LG), solid-liquid (SL), solid-gas (SG), and liquid-solid (LS). These concepts were implemented in a gradual manner through five levels of conceptual understanding (Table 1). Such method functions as the pathway of conceptual development that involves learning objectives from the lowest to the highest level of conceptual understanding (Duncan & Hmelo-Silver, 2009; Löfgren & Helldén, 2009; Hadenfeldt et al., 2013; Rogat et al., 2011). In other words, the set of levels, as mentioned previously, was adjusted to the students' needs so as to develop their conceptual

Commented [L6]:rusch model: for dicotomous data (0,1) PCM: for politomous data (categories:

1,2,3, ..

understanding. This took into account that each student might progress on different and nonlinear development of conceptual understanding; therefore, the levels, as illustrated in Table 1, was considered valid to illustrate the ideal conceptual development pathway (Neumann et al., 2013).

Table	1.	Level	of	Conceptual	Understanding	in	Explaining	Concept	of	Change	of	State	of
Matter	r			_			-	_		-			

Conceptual Understanding Level	Change of	Change of State of Matter/Item					
Conceptual Onderstanding Level	LG S	SL SG	LS				
5 Submicroscopic diagram of change	of <mark>5LG-</mark> 10	SL- 15SG	20LS				
state of matter	<mark>5</mark>	<mark>5 -5</mark>	<mark>-5</mark>				
4 Correlation between state of matter a	nd <mark>4LG-</mark> 9S	L-4 14SG	<mark>19LS</mark>				
the process of change of state of matt	er <mark>4</mark>	<mark>-4</mark>	<mark>-4</mark>				
3 Process of change of state of matter	<mark>3LG-</mark> 8S	L-3 <mark>13SG</mark>	18LS				
	<mark>3</mark>	<mark>-3</mark>	<mark>-3</mark>				
2 Concept of state of matter	<mark>2LG-</mark> 7S	L-2 12SG	17LS				
	2	<mark>-2</mark>	<mark>-2</mark>				
1 Factual phenomenon of state of matter	r <mark>1LG-</mark> 6S	L-1 <mark>11SG</mark>	16LS				
-	1	-1	-1				

Description: (LG = liquid-gas, SL = solid-liquid, SG = solid-gas, LS = liquid-gas)

Item Design and Assessment Scheme

The second phase involved an item design. In the 4TMC instrument, all the items consisted of four-tier multiple-choices. To put it another way, each item contains four questions that combine between diagnostic-summative test (Hoe & Subramaniam, 2016; Lu & Bi, 2016; Treagust, 1988) with certainty response index (hereinafter, CRI) test (Arslan et al., 2012; Hasan et al., 1999). The first-tier questions (Q1) aimed to identify whether or not the students understand the content. Moreover, questions in the second tier (Q2) were employed to clarify the students' certainty regarding their answers in the Q1. Third-tier questions (Q3) functioned to diagnose the students' reasoning regarding their answers in the Q1. Further, questions in the second tier (Q4) were employed to clarify the students' certainty regarding their answers in the Q3. Q1 and Q3 questions in each item involved five answer choices; one among them was the correct answer, while three were the distractor, and another answer choice was open-ended answer choice. This open-ended option allows the students to decide the answer by themselves, should they find no correct answer as in accordance with their conceptual

understanding. In the meantime, the Q2 and Q4 questions involved two close-ended answer choices; the first choice was for those who are uncertain of their answer, and the second choice was for the students who are very certain of their answer (Arslan et al., 2012). The distractor choices were employed in Q1 and Q3 questions to validate the diagnostic strength of the questions (Herrmann-Abell & DeBoer, 2011; Sadler, 1998). Therefore, in the Q1 and Q3 tiers, the students would have only 0.20 or 20 percent probability of choosing the correct answer.

Outcome Space and Data Collection

The third step involved the design of the outcome space of the correlation between items and construct maps (Bond & Fox, 2007; Wilson, 2009). The item validation was conducted independently by three expert validators to evaluate the extent of correlation between answer choices in Q1-Q3 in each item and the level of students' conceptual understanding. The validators were asked to clarify that the questions are easy to understand and the students' lack of linguistic competence would not hinder them from providing the right answer. The validators also required to ensure that the questions are in accordance with the syllabus, particularly with the students' conceptual understanding as based on the construct map. The questions in each item were also validated in several aspects, such as: ambiguity, time allocation, directiveness towards a particular answer, and subjective or emotional expression. Fleiss κ measure was employed to acquire information on the validators' approval. From the measure, it was generated that the κ value = 0.97, indicating that the three validators agreed that the 4TMC items were valid in correlating between the answer choices and the students' conceptual understanding.

The next step was to acquire data based on the measurement instrument. The instrument was distributed to 427 students in Gorontalo, Indonesia. The students comprised 171 (40.05%) senior high school students (or students A), 83 (19.44%) university freshmen majoring chemistry education (or students B), 66 (15.45%) second-year university students majoring

Commented [L7]: K denotes what parameter of quantity?

Commented [L8]: tested

chemistry education (or students C), 55 (12.88%) third-year university student majoring chemistry education (or students D), and 52 (12.18%) fourth-year university students majoring chemistry education (or students E). Based on gender, the female participants comprised 369 participants (86.41%), and the male counterparts consisted of 58 participants (13,58%). The participants were given no particular educational treatments and had stated their voluntary consent to participate in the research.

Rasch Model Measurement and Data Analysis

The fourth step was to conduct the Rasch model measurement. This step was implemented to define the correlation between the score generated and the students' conceptual understanding level as elaborated within the construct map. The involvement of Rasch model measurement lay on the assumption that the item difficulty level is dependent on the students' answer, and that the students' understanding is dependent on the estimation of item difficulty (Linacre, 2012).

Rasch partial credit model (PCM) was employed to evaluate the learning progress through structured questions; this took into account that the instrument items involved gradual and structured questions (Bond and Fox, 2007; Masters, 1982; Sumintono and Widhiarso, 2015; Wilson, 2009). The model was stated into the following formula: $ln[P_nik/(1-P_nik)] - B_n - D_ik$, in which P_{nik} refers to the probability of student n with B_n ability to pick correct response in the level k of item i; while D_{ik} refers to the difficulty level k of item i, or the threshold point for the test taker who scores k, not k -1. Analysis of differences based on class level of students using One-way ANOVA.

Results and Discussion

Results

Effectiveness of Measurement Instruments

Commented [L9]: how to determine effectiveness and learning progress in students

Commented [L10]: Rasch model: if answer is only true (1) and false (0) = dicotomous PCM : PCM: if the response pattern is polytomous

Commented [L11]: 1 PL, 2 PL, 3PL, 4 PL (PL=parameter logistik)

Unidimensionality is an essential indicator to evaluate the 4TMC instrument's ability to measure students' capability of explaining the concept of change of state of matter. This indicator is measured by Principal Component Analysis of the residuals to estimate the extent of variance to which the instrument is able to measure what it is supposed to measure (Sumintono & Widhiarso, 2014). he result of raw variance explained by measures of data is 38.9%, the number almost approaches the expectation value of 39.2%. The numbers indicate that the minimum unidimensionality requirements of 20% are achieved, and simultaneously, the limit of Rasch unidimension is met (approaching 40%) (Linacre, 2012; Ling Lee, Chinna, & Sumintono, 2020). Moreover, the instrument's unexplained variance values are below 7% and considered as ideal (not exceeding 15%), signifying that the item independence rate in instrument falls into "good" category.

The second step is to measure the consistency between the item difficulty level and students' conceptual understanding. The research discovers several interesting cases regarding the difference between the items and students' conceptual understanding: Firstly, there are four items identified (LG, SL, SG and LS) that measure similar constructs within each level of conceptual understanding. Despite being in the same conceptual understanding level, the items' logit is completely different. For instance, four items were discovered in level 3, each with varying logit (8SL-3 (-0.33) < 18LS-3 (-0.29) < 3LG-3 (+0.15) < 13SG-3 (+0.30)). The numbers indicate that overall, students are more capable of explaining the concept of SL state change compared to LS, LG, and SG. This condition also occurs in the level 4, in which each item has varying logit (19LS-4 (-0.37) < 9SL-4 (+0.04) = 14SG-4(+0.04) < 4LG-4 (+0.07)). Such a finding shows that the students find it easier to explain the correlation between the state of matter and the change process in LS compared to either SL, SG, or LG. Two sample cases above have illustrated that the students' conceptual understanding differs between the change process of LG (evaporation), SG (sublimation), SL (melting), and LS (freezing).

Moreover, it is found that the items in higher conceptual understanding levels tend to have lower logit than those at a lower level. As an instance, the logit of item 19SL-4 in level 4 (-0.37) is smaller than that of item 13SG-3 in level 3 (+0.30). This signifies that students find it harder to explain the item 13SG-3 compared to item 19SL-4. Thirdly, in the same concept of change of state (for example, LS), the logit of item 17LS-2 in level 2 (-0.40) is smaller than that of item 16LS-1 in level 1 (-0.16). As illustrated by the number, students find it easier to explain the SL concept in level 2 rather than to explain the concept's macroscopic fact in level 1. The findings above indicate that the students' conceptual understanding is not consistent with the item sequence. Moreover, the findings also suggest that the item difficulty level (LG, SL, SG, and LS; particularly SL (melting) and LS (freezing)) do not match the level in the construct map.

Measurement reliability

In **Rasch** analysis, the indicator of reliability is observed from the quality of students' response patterns, the instrument, and the interaction between person-item. Within this study, item separation and person separation values are employed as the indicators. The separation index is also converted to Cronbach-equivalent value with an estimation of 0-1. The summary of measurement instrument statistics is displayed in Table 2 as follows:

Table 2. Summary of fit statistics

	Student	Item
	(N=427)	(N=20)
Mean	0.26	0.00
Standard Error	0.02	0.09
Standard Deviation (SD)	0.48	0.41
Reliability	0.82	0.99
Infit mean-square	1.02	1.03
Outfit mean-square	1.05	1.05
Infit ZSTD	0.00	0.00
Outfit ZSTD	0.10	0.30
Point Raw Score to measure correlation	0.99	-0.99
Separation index (reliability)	2.10	9.54
Cronbach Alpha (KR-20): 0.84		
Data Points : 8540		
Chi-Square : 21173		
df : 8091 (p = 0.0000)		

Commented [L12]: PCM for 1PL?

How do items fit in each item based on item analysis?

From the table 2, it is generated that the total data points are 8540 with a Chi-square value of 21173 and the degree of freedom (df) of 8091 (p = 0.0000). These numbers indicate that the measurement is deemed as "very good" and "significant". The column of students and item in the table suggest whether or not the students and the item are considered fit. The average measure value of students is +0.26 logit (μ > 0.00), signifying that the students in overall are competent to explain the concept of change of state of matter. If the separation index value of students (+2.10 logit) is inputted into the person strata (H) formula, or H = [(4*separation) + 1]/3, thus, the generated H value = +3.13 (Linacre, 2012; Sumintono & Widhiarso, 2015). The person strata value (H) of 3 suggests that the students are classifiable into three groups of conceptual understanding (high, moderate, and low). On top of that, if the item's separation index value (+9.54) is processed by the same formula (H), the generated value is 13. Such a number shows that the items in the instrument are classifiable into 14 levels of difficulty. Moreover, the data illustrate that the items are deemed accurate and capable of measuring the students' competence in explaining the focused topic.

From the analysis result of students' answer pattern, the research generates Infit and Outfit MNSQ values of 1.02 and 1.05, respectively, with expectation value of 1.0. This clarifies that the students' answer pattern towards the instrument is categorized as "good". In addition, the result generates Infit ZSTD and outfit ZSTD value of 0.0 and 0.10, respectively, with an expectation value of 0.0; the numbers depict that the overall students' answer pattern is in accordance with the model. Moreover, the overall reliability of students section is 0.82, categorized as "good". From the instrument item assessment, it is generated that the Infit and Outfit MNSQ values are 1.03 and 1.05, respectively, with the expectation value of 1.0, and the Infit and Outfit ZSTD values are 0.0 and 0.3, with the expectation value of 0.0. The numbers suggest that the overall instrument is deemed as "good", proven by the instrument reliability value of 0.99. The KR-20 (alpha Cronbach) value results in 0.84, thus signifying a

good interaction between the students and the item. As acquired from the findings, the actual data in this study have met the Rasch model requirements, meaning that further analysis is considered as valid to conduct.

Level of Students' Learning Progress

The second problem of the research is: "How is the learning progress of the participants ranging from senior high school to fourth college year in explaining the focused topic?". To elaborate on that matter, the study employs data generated from the development process of 4TMC instrument to measure the students' conceptual understanding level.



Figure 1 Mean student performance level by grade

(Senior high school students = A, first-year college students = B, second-year college students = C, third-year college students = D, fourth-year college students = E)

Figure 1 displays the average competence calculated in the form of logs based on the students' academic level, ranging from A to E. The figure shows an increasing trend in students' competence development based on their respective academic level (ABCDE). Moreover, it is discovered that the group E shows better learning progress compared to the other groups (D, C, B, and A). Despite that, the One-way ANOVA test indicates a difference among the students' competence based on the academic level, in which F_{count} (6, 0142442) > F_{table} (2,39308); df = 422; p <0.05. The research, therefore, conducted a post hoc Bonferroni test to identify which group that experience significant learning progress. As extracted from the statistical result, group A and B undergo significant learning progress, while group C, D,

and E do not experience such significant advancement. This contradicts the common notion that the group CDE are college students with longer formal education experience compared to group A or B. Such finding indicates that the group CDE find it hard to explain the concept of change of state of matter.

Comparison of average competence between groups ABCDE is conducted to map out the difference in the students' learning progress in each conceptual understanding level (displayed in Table 3). The students' competence is calculated based on four items in each level of conceptual understanding. As an example, in the level 1, the students' competence is measured by referring to the mean of item 1LG-1, 6SL-1, 11SG-1, and 16SL-1; the same also applies in the next levels. Based on Table 4, it is found that the students' competence in level 1 (0.77 logit, SD = 0.86) is higher than their competence in level 2 (0.69 logit, SD = 0.95); the same also applies in the next levels. The findings above indicate that the students' conceptual understanding has not developed optimally. On top of that, the item sequence in level 1 is easier to explain compared to that in level 2. The same condition also applies in the next levels. Students find it harder to explain concepts of change of state of matter as the learning progress level increases. Simply put, the students' learning progress level is different in each level of conceptual understanding.

Table 3

Conceptual	Students' Education Level (Mean, SD)						
Understanding - Level	A (N=171)	B (N=83)	C (N=66)	D (N=55)	E (N=52)	ABDCE (N=427)	
1	0.69 (0.86)	0.80 (0.71)	0.61 (0.91)	1.29 (0.95)	1.05 (0.90)	0.77 (0.86)	
2	0.58 (1.04)	0.66 (0.75)	0.68 (0.86)	1.05 (1.00)	0.83 (0.97)	0.69 (0.95)	
3	0.19 (0.95)	0.61 (1.00)	0.33 (1.13)	0.84 (0.92)	1.10 (1.24)	0.51 (1.10)	
4	0.24 (1.00)	0.53 (0.68)	0.51 (1.12)	0.70 (0.86)	0.51 (0.71)	0.41 (0.57)	
5	-1.16	-0.80	-0.86	-0.48	-0.58	-0.84 (1.41)	

Measurement of students' average competence in each level of conceptual understanding

$(1)^{2}$ (1.40) $(1)^{1}$ $(0.0.)$ $(1)^{1}$	(1.59)	(1.46)	(1.51)	(0.85)	(1.51)
---	--------	--------	--------	--------	--------

The difference in students' learning progress levels in each conceptual understanding level depicts that longer formal education experience does not necessarily guarantee that the student will have better learning progress in explaining the focused topic. For instance, Table 4 illustrates the comparison of item logit size in level 3 that is calculated based on the students' academic level.

Table 4 Average item logit in level 3

Education		Item mean (logit) at level 3			
Level	Ν	<mark>13SG-</mark> 3	3LG-3	18LS-3	8SL-3
А	171	0.51	0.33	-0.22	-0.61
В	83	0.55	0.40	-0.43	-0.51
С	66	0.61	0.19	-0.15	-0.66
D	55	0.33	0.20	-0.15	-0.46
Е	52	0.57	0.06	-0.30	-0.33

Discussion

The result shows that: firstly, based on the logit size, the items are put in the following order: 13SG-3 > 3LG-3 > 18 SL-3 > 8SL-3. This is to say that it is harder for the students to explain the concept in item 13SG-3 compared to 3LG-3, 18SL-3, and 8SL-3. Secondly, the students' competence in each item is different and not in sequential order based on the education level (ABCDE). The finding leads to an assumption that all students in group E are supposed to perform better in explaining the item sequence in level 3 than those in group D, C, B, and A, since they progressed through longer education experience. However, the calculation result shows a different insight. In the item 13SG-3, students in group C are the most competent among all group (C (0.61) > E (0.57) > B (0.55) > A (0.51) > D (0.33)), while in the item 8SL-3, group E students are the most competent (E (-0.33) > D (-0.46) > B (-0.51) > A (-0.61) > C (-0.66)). Such a finding indicates that the students' competence is varied despite

Commented [L14]: the results of the discussion are associated with the steps to answer the research objectives (1) and (2) in this study

being at the same level. To put it another way, longer formal education tends to have an insignificant effect on the development of students' conceptual understanding.

Table 5

Category of item 13SG-3 comprehension

Grade	N	Conceptual Understanding Category - Item 13SG-3 (%)				
		LOK	<mark>AM</mark>	MFN	MFP	SK
А	171	36	21	3	20	19
В	83	19	36	8	5	31
С	66	36	27	2	8	27
D	55	13	24	4	7	53
E	52	23	12	6	12	48

Commented [L15]: Category: LOK = Lack of Knowledge, AM = All-Misconception, MFN = Misconception False Negative, MFP = Misconception False Positive, SK = Scientific Knowledge

How is the students' learning progress level in the same item? Table 5 displays the percentage data of students' competence in explaining item 13SG-3 based on five categories of conceptual understanding (LOK, AM, MFN, MFP, and SK). In the SK category, students in group D perform better among all groups (D (53%) > E (48%) > B (31%) > C (27%) >A (19%)). Simply put, more than half students in group D are capable of explaining the item 13SG-3 compared to students in other groups. Meanwhile, in LOK, students in group A and C show higher percentage among all groups (A (36%) = C (36%) > E (23%) > B (19%) > D(13%)). In other words, more than one-third of students in group A or C is incapable of explaining the item 13SG-3 compared to students in other groups due to the limited knowledge on the item. Moroever, in AM, group B shows highest percentage among all groups (B (36%) > C (27%) > D (24%) > A (21%) > E (12%)); it signifies that more than one-fourth of students in group B are incapable of explaining item 13SG-3 compared to other groups due to the misconception on the item. Such findings indicate that the high percentage in LOK and AM category is seen as one of the reasons why the students' competence is different in explaining the same item 13SG-3. To put it another way, the students' learning progress does not develop optimally in explaining item 13SG-3 due to lack of knowledge (LOK) or misconception (AM) on the item.

Commented [L16]: add an example of the problem that is discussed



Figure 2(a) Probability Category Curve of item 13SG-3 of group A, and Figure 2(b) Probability Category Curve of item 13SG-3 of group D (Category: LOK = Lack of Knowledge, AM = All-Misconception, MFN = Misconception False Negative, MFP = Misconception False Positive, SK = Scientific Knowledge)

Figure 2 illustrates the comparison of the probability category curve (PCC) of students in group A and D in item 13SG-3. The five curve shapes are the visual representation of the distribution of five categories of students' conceptual understanding. From the curves, one can identify which groups that tend to show LOK and AM category traits. It is worth noting that the curve 2(a) and 2(b) tend to be different based on the MFP curve shape, while others are relatively similar. The MFP curve of students A has a higher probability compared to that of students D; simply put, a senior high school student tends to show stronger MFP category compared to a third-year college student. The notion is supported by the finding that senior high school students are relatively incapable of providing correct reason on item 13SG-3 compared to third-year college students. On the other hand, students with low ability in group D tend to show similar curve shape of LOK, AM, and MFN with group A. This implies that both groups' conceptual understanding in the item is relatively similar. In other words, the learning progress of group D, particularly in students with low ability, has not developed
optimally despite the fact that that group D consists of third-year college students that progressed through three years of formal education experience in university.

This echoes previous findings that the learning progress is highly dependent on the students' learning process and experience (Duschl et al., 2011; Park et al., 2017; Wilson, 2009). Learning progress is defined as a sophisticated and systematic way of thinking, in which the students will undergo gradual progress when learning a topic for a long time interval. Such a systematic way of thinking is formed by the learning practices and education experience (Emden et al., 2018). On top of that, the research findings are in line with previous studies that highlighted that students have distinctive comprehension formed by their own experience (Chi et al., 2018; Emden et al., 2018; Hoe & Subramaniam, 2016; Jin, Mikeska, Hokayem, & Mavronikolas, 2019; Rogat et al., 2011; Testa et al., 2019). Such distinctive knowledge has not been explored by evaluation or intervention through learning roadmaps that are in accordance with remedial learning (Smith et al., 2006). In spite of that, it is considered essential to conduct a further analysis that focuses on the modification of conceptual understanding category and analysis variation that is able to define the characteristics of students' alternative conception. The development procedures, as explained in the methodology, has resulted in 4TMC instrument; however, instrument development is seen as an essential continuous process (Wilson, 2009, 2012).

Based on the research findings, the study identifies several important notes on the development of the 4TMC instrument. Firstly, further analysis of the characteristic of students' response behavior is necessary to conduct regarding the item clarity and the measured concept. The findings have implied that the percentage of LOK and AM understanding category is relatively dominant and tends to increase along with the level of conceptual understanding. Hence, the development of the concept level requires taking into consideration any potential term use that might confuse the students. A further study on the identification of commonly-understood terms or concepts is therefore essential. Secondly, a

separate analysis is required to diagnose the factors contributing to the students' lack of knowledge and misconception. Regarding that, further analysis can be conducted by applying the analysis methods developed by previous studies (Caleon & Subramaniam, 2010; Hoe & Subramaniam, 2016; Oon & Subramaniam, 2013). Thirdly, it is discovered that the concepts LG, SG, SL and LS were interpreted differently by the students. Despite being in the same conceptual understanding level, the items' difficulty level are completely different. Therefore, an evaluation on answer choices requires one to focus on the representation of understanding at the same level.

One of the features of the Rasch model is that the model facilitates one to identify any correlation between the construct map and the students' competence in ways that the students' competence can be analyzed by referring to the difference in item difficulty level. The 4TMC instrument indicates that there are students with very high ability as well as students with low ability in each group. Such a gap serves as the basis for qualitative interpretation to elaborate on the difference in students' competence. The insight is applicable in the learning process of chemistry subject. The instrument is expected to be beneficial for teachers in developing a formative test to identify the students' progress of conceptual understanding. On top of that, teachers are able to implement the instrument as a diagnostic instrument to evaluate students' conceptual understanding in providing feedback on their learning progress. Further, the teachers will be able to develop instructional strategies that are specifically designed to tackle the students' difficulty in developing an epistemological explanation regarding the concept of change of state of matter.

Conclusions

The result revealed that the integration of the 4TMC test and Rasch modeling is effective and valid to be treated as the diagnostic instrument to measure students' learning progress. Moreover, it is discovered that students in group A, B, C, D, and E, particularly those with

Commented [L17]: conclusions adjusted to the research objectives

low ability, are hampered in developing an epistemological explanation of the concept. This blames the students' lack of certainty in their answer and reason; thus, assumed as having lack of knowledge or misconception. The low-ability students' curve shape of LOK and AM is consistent in the competence interval of less than 0.1 logit. On the other hand, the students' ability gets lower as the conceptual understanding level increases. Such finding indicates that the learning process and education experience provide a limited contribution for the students in developing a systematic way of thinking regarding the concept of change of state of matter.

Recommendations

The Based on the results of the study, there are several recommendations for researchers and teachers. For researchers, the findings of this research can be followed up to examine more in how students build their understanding gradually in explaining the concept of particles in substance form changes. The study can be conducted by developing tests that aim to evaluate and diagnose the process of student knowledge formation and development while being able to identify at the level of education where the confusion of understanding occurs. The evaluation becomes more objective, not only reviewed from the student's point of ability but can be reviewed from the teacher's ability. The model of *Rasch's* multi-faced item response pattern approach becomes one of the important parts recommended for such objectives. In this way, students' ability to develop epistemological knowledge, and their ability to significantly actualize the knowledge gained can be measured well.

On the other hand, for teachers, the results of this study along with the stages of analysis approach used can be a reference in evaluating the progress of learners' learning, as well as determining alternative thinking frameworks of students in explaining the concept of substance change. The information serves as strategic feedback in formulating instructional strategies and preparing remedial learning, especially for students who have difficulty in developing epistemological explanations of substance changes.

Limitations

The limitations of the research are primarily related to the misrepresentation of student reasoning, which may arise in its efforts to connect phenomena and concepts measured in each item. In this context, the student may not excel to explain, because of his incapableness in using his heuristic reasoning. This instrument is not equipped with items that evaluate the heuristic abilities of the student in question. However, researchers decided to record this incompetence as a misconception or vague knowledge. For further research, it is recommended that the instrument be equipped with items that measure students' emotional and heuristic reasoning according to the conceptual framework to be evaluated.

Acknowledgments

The researchers would like to express their gratitude towards the Directorate of Research and Community Service, Ministry of Research and Technology of Republic of Indonesia, for the financial support through the University Basic Research Excellence Grant Program in the Research and Community Service Office of Universitas Negeri Gorontalo, 2020.

References

- Aktan, D. C. (2013). Investigation of students' intermediate conceptual understanding levels: The case of direct current electricity concepts. *European Journal of Physics*, 34(1), 33– 43. https://doi.org/10.1088/0143-0807/34/1/33
- Arslan, H. O., Cigdemoglu, C., & Moseley, C. (2012). A three-tier diagnostic test to assess pre-service teachers' misconceptions about global warming, greenhouse effect, ozone layer depletion, and acid rain. *International Journal of Science Education*, 34(11), 1667– 1686. https://doi.org/10.1080/09500693.2012.680618
- Bond, T. G., & Fox, C. M. (2007). *Applying The Rasch Model: Fundamental Measurent in the Human Sciences* (2nd Ed.). Routledge Taylor & Francis Group

- Caleon, I. S., & Subramaniam, R. (2010). Do students know what they know and what they don't know? Using a four-tier diagnostic test to assess the nature of students' alternative conceptions. *Research in Science Education*, 40(3), 313–337. https://doi.org/10.1007/s11165-009-9122-4
- Chandrasegaran, A. L., Treagust, D. F., & Mocerino, M. (2007). The development of two tier multiple-choice diagnostic instrument for evaluating secondary school students' ability to describe and explain chemical reactions using multiple levels of representation. *Chemistry Education Research and Practice*, 8(3), 293–307
- Chi, S., Wang, Z., Luo, M., Yang, Y., & Huang, M. (2018). Student progression on chemical symbol representation abilities at different grade levels (Grades 10–12) across gender. *Chemistry Education Research and Practice*, 19(4), 1055–1064. https://doi.org/10.1039/c8rp00010g
- Claesgens, J., Scalise, K., Wilson, M., & Stacy, A. (2009). Mapping student understanding in chemistry: The perspectives of chemists. *Science Education*, 93(1), 56–85. https://doi.org/10.1002/sce.20292.
- Duncan, R. G., & Hmelo-Silver, C. E. (2009). Learning progressions: Aligning curriculum, instruction, and assessment. *Journal of Research in Science Teaching*, 46(6), 606–609. https://doi.org/10.1002/tea.20316
- Duschl, R., Maeng, S., & Sezen, A. (2011). Learning progressions and teaching sequences: A review and analysis. *Studies in Science Education*, 47(2), 123–182. https://doi.org/10.1080/03057267.2011.604476
- Emden, M., Weber, K., & Sumfleth, E. (2018). Evaluating a learning progression on "Transformation of Matter" on the lower secondary level. *Chemistry Education Research and Practice*, 19(4), 1096–1116. https://doi.org/10.1039/c8rp00137e
- Habiddin, & Page, E. M. (2019). Development and validation of a four-tier diagnostic instrument for chemical kinetics (FTDICK). *Indonesian Journal of Chemistry*, 19(3), 720–736. https://doi.org/10.22146/ijc.39218
- Hadenfeldt, J. C., Bernholt, S., Liu, X., Neumann, K., & Parchmann, I. (2013). Using ordered multiple-choice items to assess students' understanding of the structure and composition of matter. *Journal of Chemical Education*, 90(12), 1602–1608. https://doi.org/10.1021/ed3006192
- Hasan, S., Bagayoko, D., & Kelley, E. L. (1999). Misconceptions and the certainty of response index (CRI). *Physics Education*, 34(5), 294–299. https://doi.org/10.1088/0031-9120/34/5/304
- Herrmann-Abell, C. F., & Deboer, G. E. (2016). Using rasch modeling and option probability curves to diagnose students' misconceptions. *American Eduacational Research Assossiation*, 1–12
- Hoe, K. Y., & Subramaniam, R. (2016). On the prevalence of alternative conceptions on acidbase chemistry among secondary students: Insights from cognitive and confidence measures. *Chemistry Education Research and Practice*, 17(2), 263–282.

https://doi.org/10.1039/c5rp00146c

- Jin, H., Mikeska, J. N., Hokayem, H., & Mavronikolas, E. (2019). Toward coherence in curriculum, instruction, and assessment: A review of learning progression literature. *Science Education*, 103(5), 1206–1234. https://doi.org/10.1002/sce.21525
- Klassen, S. (2006). Contextual assessment in science education: Background, issues, and policy. *Science Education*, 90(5), 820–851. https://doi.org/10.1002/sce.20150
- Laliyo, Botutihe, & Panigoro. (2019). The development of two-tier instrument based on distractor to assess conceptual understanding level and student misconceptions in explaining redox reactions. *International Journal of Learning, Teaching and Educational Research*, 18(9), 216–237. https://doi.org/10.26803/ijlter.18.9.12
- Linacre, J. M. (2012). A user's guide to WINSTEPS ® MINISTEP Rasch-model computer program: Program manual 3.75.0. https://doi.org/ISBN 0-941938-03-4
- Linacre, J. M. (2020). A User's Guide to WINSTEPS ® MINISTEP Rasch-Model Computer Programs Program Manual 4.5.1. https://doi.org/ISBN 0-941938-03-4
- Ling Lee, W., Chinna, K., & Sumintono, B. (2020). Psychometrics assessment of HeartQoL questionnaire: A Rasch analysis. *European Journal of Preventive Cardiology*, 1–6. https://doi.org/10.1177/2047487320902322
- Liu, X. (2012). Developing measurement instruments for science education research. In B. Fraser, K. G. Tobin, & C. J. McRobbie (Eds.), Second international handbook of science education (pp. 651–665). Springer Netherlands
- Lu, S., & Bi, H. (2016). Development of a measurement instrument to assess students' electrolyte conceptual understanding. *Chemistry Education Research and Practice*, 17(4), 1030–1040. https://doi.org/10.1039/c6rp00137h
- Morell, L., Collier, T., Black, P., & Wilson, M. (2017). A construct-modeling approach to develop a learning progression of how students understand the structure of matter. *Journal of Research in Science Teaching*, 54(8), 1024–1048. https://doi.org/10.1002/tea.21397
- Neumann, K., Viering, T., Boone, W. J., & Fischer, H. E. (2013). Towards a learning progression of energy. *Journal of Research in Science Teaching*, 50(2), 162–188. https://doi.org/10.1002/tea.21061
- Park, M., Liu, X., & Waight, N. (2017). Development of the connected chemistry as formative assessment pedagogy for high school chemistry teaching. *Journal of Chemical Education*, 94(3), 273–281. https://doi.org/10.1021/acs.jchemed.6b00299
- Peterson, R. F., Treagust, D. F., & Garnett, P. (1989). Development and application of a diagnostic instrument to evaluate grade-11 and -12 students' concepts of covalent bonding and structure following a course of instruction. *Journal of Research in Science Teaching*, 26(4), 301–314. https://doi.org/10.1002/tea.3660260404
- Rogat, A., Anderson, C., Foster, J., Goldberg, F., Hicks, J., Kanter, D., ... Wiser, M. (2011). Developing learning progression in support of the new science standards: A RAPID

workshop series. (4), 163. https://doi.org/10.12698/cpre.2011.lprapid

- Sadler, P. M. (1999). The relevance of multiple-choice testing in assessing science understanding. In J. J. Mintzes, J. H. Wandersee, & J. D. Novak (Eds.), Assessing science understanding: A human constructivist view (pp. 251–274). Zaccheus Onumba Dibiaezue Memorial Libraries. https://zodml.org/sites/default/files/%5BJoel_J._Mintzes%2C_James_H._Wandersee%2 C_Joseph_D._No_0.pdf
- Smith, C. L., Wiser, M., Anderson, C. W., & Krajcik, J. (2006). Implications of research on children's learning for standards and assessment: A proposed learning progression for matter and the atomic-molecular theory. *Measurement: Interdisciplinary Research & Perspective*, 4(1–2), 1–98. https://doi.org/10.1080/15366367.2006.9678570
- Sumintono, B., & Widhiarso, W. (2014). Aplikasi model Rasch untuk penelitian ilmu-ilmu sosial [Application of Rasch model in social science research]. Trim Komunikata. https://www.researchgate.net/publication/268688933%0AAplikasi
- Testa, I., Capasso, G., Colantonio, A., Galano, S., Marzoli, I., Scotti di Uccio, U., ... Zappia, A. (2019). Development and validation of a university students' progression in learning quantum mechanics through exploratory factor analysis and Rasch analysis. *International Journal of Science Education*, 41(3), 388–417. https://doi.org/10.1080/09500693.2018.1556414
- Treagust, D. F. (1988). Development and use of diagnostic tests to evaluate students' misconceptions in science. *International Journal of Science Education*, 10(2), 159–169. https://doi.org/10.1080/0950069880100204
- Tyson, L., Treagust, D. F., & Bucat, R. B. (1999). The complexity of teaching and learning chemical equilibrium. *Journal of Chemical Education*, 76(2–4), 554–558. https://doi.org/10.1021/ed077p1560.1
- Wilson, M. (2005). Constructing measures: an item response modeling approach. Lawrence Erlbaum Associates, Inc. https://doi.org/10.4324/9781410611697
- Wilson, M. (2008). Cognitive diagnosis using item response models. Zeitschrift Für Psychologie / Journal of Psychology, 216(2), 74–88. https://doi.org/10.1027/0044-3409.216.2.74
- Wilson, M. (2009). Measuring progressions: Assessment structures underlying a learning progression. Journal of Research in Science Teaching, 46(6), 716–730. https://doi.org/10.1002/tea.20318
- Wilson, M. (2012). Responding to a challenge that learning progressions pose to measurement practice. In A. C. Alonzo & A. W. Gotwals (Eds.), *Learning progression* in science (pp. 317–344). https://doi.org/10.1007/978-94-6091-824-7

Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress

Abstract: The purpose of this study was to developing and implementation four-tier multiple-choice (hereinafter, 4TMC) instrument with Rasch model to evaluate students' learning progress in explaining the concept of change of state of matter. [The data were obtained through development and validation techniques on 20 4TMC items distributed to 427 students]. On each item, the study applied diagnostic-summative assessment and certainty response index. The students' conceptual understanding level was categorized based on the combination their answer choices; the measurement generated Partial-Credit polytomous Rasch model data. The data were further processed by **WINSTEPS version 4.5.3** software to equate the data interval rate. Analysis of differences based on class level of students using Analysis of Variants (One-way ANOVA). The result revealed that the integration of 4TMC test and Rasch modeling was effective to be treated as the instrument to measure students' learning progress. One-way ANOVA test indicates a difference among the students' competence based on the academic level, in which F_{count} (6, 0142442) > F_{table} (2,39308); df = 422; p <0.05. On top of that, it was discovered that low-ability students see very slow progress due to the lack of knowledge as well as a misconception in explaining the concept as mentioned above.]

Keywords: Learning progress. four-tier, change of state of matter, Rasch model.

abstract. What do you readers want to know? So, adjust the order with the following Abstract writing style steps: 1. Why did the study need to be done? - Introduce topic and problem shortly. It's important to state the background before

Commented [A1]: Give good first impression from your

the research aims.

2. What did you do? - your aims and methodology.

3. What did you find? - key results

4. How study will advance the field? – conclusions and implications.

Commented [A2]: The research design is not yet visible

Commented [A3]: The name of the software for analyzing data should not be written in the abstract but described in the methodology section.

Introduction

Central to the notion of science learning is the development of students' scientific understanding of basic concepts of sciences (Hadenfeldt et al., 2013), particularly, change of state of matter (Emden et al., 2018). Aside from the issue, several studies have also highlighted the students' inability to provide an epistemological explanation of basic concepts of sciences (Chi et al., 2018). Efforts to solve the issues, however, have shown little progress, as the students might have more complex perceptions regarding the alternative concept they understand (Morell et al., 2017).

Education practitioners have recommended the utilization of learning progress concept as the instructional method to provide guidance and direction and to adjust the curriculum with the learning process and assessment (Claesgens *et al.*, 2009; Duncan & Hmelo-Silver, 2009;

Commented [A4]: What concepts are described on the instrument?

Commented [A5]: At the end of the Abstract, you must add at least one remarkable recommendation.

Commented [A6]: For suggestions and please note: 1.The introduction does not need to be wordy, in paragraph 1 make the urgency of the misconceptions that have occurred lately

2.In the next paragraph, come up with the state of the art. There have been so many studies that have resulted in misconception instruments, come up with those studies.

3.Next, compare the results of your research (instrument) with the results of previous studies. Besides the material, what is the difference between your research and previous research?

4. What is the specialty of your instrument compared to other studies? Why should people choose your instrument over other people's development instruments?

Rogat *et al.*, 2011). Learning progress is defined as a sophisticated and systematic way of thinking. This method is applicable for a learning process, in which the students will undergo gradual progress when learning a topic in a long duration. Its effectiveness is highly dependent on the learning process and the students' learning experience (Duschl et al., 2011). The concept involves certain sets of gradual levels that represent conceptual understanding, ranging from low level up to comprehensive level.

The notion of learning progress is highly distinctive to each student and is dependent to one's learning experience (Rogat et al., 2011); therefore, there is no learning roadmap that is suitable for all kinds of students (Smith et al., 2006). Each student constructs one's understanding in a different way; moreover, the construction process is varied depending on the students' conceptual understanding level (Aktan, 2013). This is to say that each student undergoes a different rate of learning progress, understanding level, and knowledge construction. Simply put, the development of scientific comprehension among students is not linear (Neumann et al., 2013). Therefore, this study regards each level of students' conceptual understanding for more advanced level of understanding (Hadenfeldt et al., 2013). A student who faces difficulty in a certain level of understanding will see a lack of progress to a more advanced level. This in turn hinders the student's ability to construct an epistemological explanation on the basic concepts of science.

Efforts to diagnose the epistemological problems, as mentioned previously, are feasible to conduct if the extent of students' conceptual understanding is formulated. Within this context, the learning progress is treated as the method to evaluate students' conceptual understanding. The diagnostic information generated is reliable to be treated as a reference for the teachers in developing accurate and valid instructional components to guide the students to progress to the next level. Despite the potentials, this study deems that it is challenging for the teachers to construct such an accurate instrument.

Commented [A7]:

- Consider summarising the text based on the study purpose.
 Focus more on the empirical studies' backgrounds or studies' problem.
- Add more information to enable readers' understanding of the authors' view.
- You are encouraged to write concisely. The text can be reduced significantly.

Among the diagnostic instruments that are considered applicable is the 4TMC instrument. It is the development of two-tier multiple-choice test recommended by Treagust (1988) and Chandrasegaran et al., (2007). The use of two-tier instrument is familiar in identifying students' understanding in select topics such as electrochemistry (Lu & Bi, 2016), covalent bond (Peterson, Treagust, & Garnett, 1989), and chemical equilibrium (Tyson et al., 1999). Despite its reputation in academia, the two-tier test has raised criticism due to its sole focus on the facts and negligence towards students' understanding (Klassen, 2006). Therefore, several experts propose the renewed version of the test by adding distractor answer choices to strengthen the diagnostic value of the items (Herrmann-Abell & DeBoer, 2011; Sadler, 1998). In addition, some have highlighted the test's weakness in cases where students' tended to pick the answer choice and the reasoning randomly. This illustrates that the students were uncertain and possessed several misconceptions in the first tier question. In such cases, teachers faced difficulty in differentiating between guessed answers and misconceptions (Habiddin & Page, 2019; Hasan et al., 1999).

The criticism laid against the model has sparked the innovation of three-tier and four-tiers instruments. Both instruments feature two multi-level questions, also similar with two-tier test. In the three-tier test, however, the measurement of students' certainty level is conducted simultaneously in both first and second-tier questions; in the meantime, the measurement is conducted separately in the first two tiers (Caleon & Subramaniam, 2010). The value of students' certainty rate ranges from one (very uncertain) to five (very certain).

Three-tier test lacks validity in measuring the students' certainty rate regarding both the answer choice and the reasoning, whether or not the value of certainty rate refers only to the answer choice, to the reasoning, or both. Such weakness will in turn obstructs the evaluation and classification process of students' responses (Arslan et al., 2012). In the four-tier instrument, the measurement of certainty rate also involves the answer choice in the first tier and the reasoning in the third tier (Arslan et al., 2012; Loh et al., 2014). Regarding this

feature, four-tier test is considered more accurate than the three-tier test. Students who pick wrong answer choices with high certainty indicate that they have a very high misconception on the measured item (Hoe & Subramaniam, 2016).

Four-tier instruments are used in studies discussing topics such as physics education (Caleon & Subramaniam, 2010), chemical thermodynamics (Sreenivasulu & Subramaniam, 2013), transition metal (Sreenivasulu & Subramaniam, 2014), acid-base reaction (Hoe & Subramaniam, 2016), and chemical kinetics (Habiddin & Page, 2019). However, it is worth noticing that studies on chemistry topic which employ four-tiers instruments tend to focus on describing alternative conception. To put it another way, the higher the certainty rate is, the stronger the students' alternative conception will be. Despite its potentials, the scholarly discussion has overlooked the implementation of a four-tier diagnostic instrument to measure students' learning progress. Therefore, further analysis is essential on the application of 4TMC test in several domains analyzes by Rasch model approach.

The use of Rasch model has been introduced since the 2000s in the science education research; it features the instrument that integrates diagnostic assessment and summative assessment (Liu, 2012; Wei et al., 2012). On top of that, the diagnostic assessment approach is introduced to conduct an in-depth analysis of the construction process of students' conceptual understanding (Claesgens et al., 2009; Hadenfeldt et al., 2013; Lu & Bi, 2016). [This study employs 4MTC test and Rasch modeling as a diagnostic tool to evaluate students' learning progress in explaining the change of state of matter]. The study focus revolves around two research questions: [1) How is the effectiveness of 4TMC instrument to evaluate the students' learning progress in explaining concepts of change of state of matter. 2) How is the learning progress in students ranging from the senior high school level up to the senior (fourth) year of college in explaining the concepts?]

Commented [A8]: Is this novelty of your research just limited to matter?

Commented [A9]: The abstract section states, "The purpose of this study is to develop ..." Why is the purpose of the Abstract not compatible with the study question? Attention to the focus of your study purposes and synchronize it with the 4TMC instrument you developed. Where is the Rasch model involved?

Methodology

Development Model

This research used a development research referring to the test development model from Wilson. Wilson (2005, 2008) introduces four steps of measurement instrument development: The first step is to the learning progress variable focused on a characteristic measured at a particular time unit. The second step comprises the design process of items or tasks used to measure students' responses. Moreover, the third step involves outcome space, in which the students' responses are categorized into all items related with the learning progress variable. On top of that, the fourth step employs measurement model, such as Rasch model. This recommendation is proven valid to be implemented in developing measurement instrument for different construct variables (Barbera, 2013; Chi et al., 2018; Hadenfeldt et al., 2013; Laliyo, Botutihe, & Panigoro, 2019; Lu & Bi, 2016; Wei et al., 2012; Wilson, 2009; Wind, Tsai, Grajeda, & Bergin, 2018). The study conducted development of measurement instrument instruments by referring to Wilson's recommendation (2005, 2008) and adopted Treagust's framework (1988) of item development. The present study also included two questions related to certainty rate (Arslan et al., 2012; Habiddin & Page, 2019: Hasan et al., 1999). The obtained data were analyzed by Rasch model approach.

Construct Map: Determining Level of Understanding

The first step was to develop the construct of measured variables. The study involved four concepts of change of state of matter: liquid-gas (LG), solid-liquid (SL), solid-gas (SG), and liquid-solid (LS). These concepts were implemented in a gradual manner through five levels of conceptual understanding (Table 1). Such method functions as the pathway of conceptual development that involves learning objectives from the lowest to the highest level of conceptual understanding (Duncan & Hmelo-Silver, 2009; Löfgren & Helldén, 2009; Hadenfeldt et al., 2013; Rogat et al., 2011). In other words, the set of levels, as mentioned previously, was adjusted to the students' needs so as to develop their conceptual

Commented [A10]:

The explanation should be given in the form of a research/development procedure flowchart. This way will make it easier for readers or other researchers to understand the flowchart of your study.

Commented [A11]: This statement is a repetition of a sentence like the beginning of a paragraph in the development model.

Commented [A12]: What software for the Rasch model used, and what is the version?

Include the name of the software and the version of the software used to analyze the data.

understanding. This took into account that each student might progress on different and nonlinear development of conceptual understanding; therefore, the levels, as illustrated in Table 1, was considered valid to illustrate the ideal conceptual development pathway (Neumann et al., 2013).

Table 1. Level of Conceptual Understanding in Explaining Concept of Change of State of Matter

Conceptual Understanding Level –		Change of State of Matter/Item				
		LG	SL	SG	LS	
5	Submicroscopic diagram of change of	5LG-	10SL-	15SG	20LS	
	state of matter	5	5	-5	-5	
4	Correlation between state of matter and	4LG-	9SL-4	14SG	19LS	
	the process of change of state of matter	4		-4	-4	
3	Process of change of state of matter	3LG-	8SL-3	13SG	18LS	
		3		-3	-3	
2	Concept of state of matter	2LG-	7SL-2	12SG	17LS	
		2		-2	-2	
1	Factual phenomenon of state of matter	1LG-	6SL-1	11SG	16LS	
	-	1		-1	-1	

Description: (LG = liquid-gas, SL = solid-liquid, SG = solid-gas, LS = liquid-gas)

Item Design and Assessment Scheme

The second phase involved an item design. In the 4TMC instrument, all the items consisted of four-tier multiple-choices. To put it another way, each item contains four questions that combine between diagnostic-summative test (Hoe & Subramaniam, 2016; Lu & Bi, 2016; Treagust, 1988) with certainty response index (hereinafter, CRI) test (Arslan et al., 2012; Hasan et al., 1999). [The first-tier questions (Q1) aimed to identify whether or not the students understand the content. Moreover, questions in the second tier (Q2) were employed to clarify the students' certainty regarding their answers in the Q1. Third-tier questions (Q3) functioned to diagnose the students' reasoning regarding their answers in the Q1. Further, questions in the second tier (Q4) were employed to clarify the students' certainty regarding in each item involved five answer choices; one among them was the correct answer, while three were the distractor, and another answer choice was open-ended answer choice. This open-ended option allows the students to decide the answer by themselves, should they find no correct answer as in accordance with their conceptual

understanding. In the meantime, the Q2 and Q4 questions involved two close-ended answer choices; the first choice was for those who are uncertain of their answer, and the second choice was for the students who are very certain of their answer (Arslan et al., 2012). The distractor choices were employed in Q1 and Q3 questions to validate the diagnostic strength of the questions (Herrmann-Abell & DeBoer, 2011; Sadler, 1998). Therefore, in the Q1 and Q3 tiers, the students would have only 0.20 or 20 percent probability of choosing the correct answer.]

Outcome Space and Data Collection

The third step involved the design of the outcome space of the correlation between items and construct maps (Bond & Fox, 2007; Wilson, 2009). The item validation was conducted independently by three expert validators to evaluate the extent of correlation between answer choices in Q1-Q3 in each item and the level of students' conceptual understanding. The validators were asked to clarify that the questions are easy to understand and the students' lack of linguistic competence would not hinder them from providing the right answer. The validators also required to ensure that the questions are in accordance with the syllabus, particularly with the students' conceptual understanding as based on the construct map. The questions in each item were also validated in several aspects, such as: ambiguity, time allocation, directiveness towards a particular answer, and subjective or emotional expression. Fleiss κ measure was employed to acquire information on the validators' approval. From the measure, it was generated that the κ value = 0.97, indicating that the three validators agreed that the 4TMC items were valid in correlating between the answer choices and the students' conceptual understanding.

The next step was to acquire data based on the measurement instrument. The instrument was distributed to 427 students in Gorontalo, Indonesia. The students comprised 171 (40.05%) senior high school students (or students A), 83 (19.44%) university freshmen majoring chemistry education (or students B), 66 (15.45%) second-year university students majoring

Commented [A13]: Change this explanation in the form of a picture or an explanation in the form of a table, so that it is shorter. There have been many studies using four-tier instruments, showing differences as a form of Novelty.

Next, attach the 4TMC diagnostic test instrument that has been developed in the Appendix session at the end of the paper (after Reference) for the sake of transparency in the development and form of **Novelty** of this paper.

Commented [A14]: The research design is not yet visible

Commented [A15]: Write down the sampling technique used in this study.

chemistry education (or students C), 55 (12.88%) third-year university student majoring chemistry education (or students D), and 52 (12.18%) fourth-year university students majoring chemistry education (or students E). Based on gender, the female participants comprised 369 participants (86.41%), and the male counterparts consisted of 58 participants (13,58%). The participants were given no particular educational treatments and had stated their voluntary consent to participate in the research.

Rasch Model Measurement and Data Analysis

The fourth step was to conduct the Rasch model measurement. This step was implemented to define the correlation between the score generated and the students' conceptual understanding level as elaborated within the construct map. The involvement of Rasch model measurement lay on the assumption that the item difficulty level is dependent on the students' answer, and that the students' understanding is dependent on the estimation of item difficulty (Linacre, 2012).

Rasch partial credit model (PCM) was employed to evaluate the learning progress through structured questions; this took into account that the instrument items involved gradual and structured questions (Bond and Fox, 2007; Masters, 1982; Sumintono and Widhiarso, 2015; Wilson, 2009). The model was stated into the following formula: $ln[P_nik/(1-P_nik)] - B_n - D_ik$, in which P_{nik} refers to the probability of student n with Bn ability to pick correct response in the level k of item i; while Dik refers to the difficulty level k of item i, or the threshold point for the test taker who scores k, not k -1. Analysis of differences based on class level of students using One-way ANOVA.

Results and Discussion

Results

Effectiveness of Measurement Instruments

Commented [A16]: Does gender influence research results? Explain the findings on the research results

Commented [A17]: You need to show examples of the instruments you've made.

After that an analysis of these diagnostic questions should also be included. Why can this question (4TMC) be called a misconception measurement instrument? What are the criteria for the 4TMC test instrument on Bloom's taxonomy? Explain in the discussion section Unidimensionality is an essential indicator to evaluate the 4TMC instrument's ability to measure students' capability of explaining the concept of change of state of matter. This indicator is measured by Principal Component Analysis of the residuals to estimate the extent of variance to which the instrument is able to measure what it is supposed to measure (Sumintono & Widhiarso, 2014). he result of raw variance explained by measures of data is 38.9%, the number almost approaches the expectation value of 39.2%. The numbers indicate that the minimum unidimensionality requirements of 20% are achieved, and simultaneously, the limit of Rasch unidimension is met (approaching 40%) (Linacre, 2012; Ling Lee, Chinna, & Sumintono, 2020). Moreover, the instrument's unexplained variance values are below 7% and considered as ideal (not exceeding 15%), signifying that the item independence rate in instrument falls into "good" category.

The second step is to measure the consistency between the item difficulty level and students' conceptual understanding. The research discovers several interesting cases regarding the difference between the items and students' conceptual understanding: Firstly, there are four items identified (LG, SL, SG and LS) that measure similar constructs within each level of conceptual understanding. Despite being in the same conceptual understanding level, the items' logit is completely different. For instance, four items were discovered in level 3, each with varying logit (8SL-3 (-0.33) < 18LS-3 (-0.29) < 3LG-3 (+0.15) < 13SG-3 (+0.30)). The numbers indicate that overall, students are more capable of explaining the concept of SL state change compared to LS, LG, and SG. This condition also occurs in the level 4, in which each item has varying logit (19LS-4 (-0.37) < 9SL-4 (+0.04) = 14SG-4(+0.04) < 4LG-4 (+0.07)). Such a finding shows that the students find it easier to explain the correlation between the state of matter and the change process in LS compared to either SL, SG, or LG. Two sample cases above have illustrated that the students' conceptual understanding differs between the change process of LG (evaporation), SG (sublimation), SL (melting), and LS (freezing).

Moreover, it is found that the items in higher conceptual understanding levels tend to have lower logit than those at a lower level. As an instance, the logit of item 19SL-4 in level 4 (-0.37) is smaller than that of item 13SG-3 in level 3 (+0.30). This signifies that students find it harder to explain the item 13SG-3 compared to item 19SL-4. Thirdly, in the same concept of change of state (for example, LS), the logit of item 17LS-2 in level 2 (-0.40) is smaller than that of item 16LS-1 in level 1 (-0.16). As illustrated by the number, students find it easier to explain the SL concept in level 2 rather than to explain the concept's macroscopic fact in level 1. The findings above indicate that the students' conceptual understanding is not consistent with the item sequence. Moreover, the findings also suggest that the item difficulty level (LG, SL, SG, and LS; particularly SL (melting) and LS (freezing)) do not match the level in the construct map.

Measurement reliability

In Rasch analysis, the indicator of reliability is observed from the quality of students' response patterns, the instrument, and the interaction between person-item. Within this study, item separation and person separation values are employed as the indicators. The separation index is also converted to Cronbach-equivalent value with an estimation of 0-1. The summary of measurement instrument statistics is displayed in Table 2 as follows:

Table 2. Summary of fit statistics

	Student	Item
	(N=427)	(N=20)
Mean	0.26	0.00
Standard Error	0.02	0.09
Standard Deviation (SD)	0.48	0.41
Reliability	0.82	0.99
Infit mean-square	1.02	1.03
Outfit mean-square	1.05	1.05
Infit ZSTD	0.00	0.00
Outfit ZSTD	0.10	0.30
Point Raw Score to measure correlation	0.99	-0.99
Separation index (reliability)	2.10	9.54
Cronbach Alpha (KR-20): 0.84		
Data Points : 8540		
Chi-Square : 21173		
df : 8091 (p = 0.0000)		

Commented [A18]: delete these words

From the Table 2, it is generated that the total data points are 8540 with a Chi-square value of 21173 and the degree of freedom (df) of 8091 (p = 0.0000). These numbers indicate that the measurement is deemed as "very good" and "significant". The column of students and item in the table suggest whether or not the students and the item are considered fit. The average measure value of students is +0.26 logit (μ > 0.00), signifying that the students in overall are competent to explain the concept of change of state of matter. If the separation index value of students (+2.10 logit) is inputted into the person strata (H) formula, or H = [(4*separation) + 1]/3, thus, the generated H value = +3.13 (Linacre, 2012; Sumintono & Widhiarso, 2015). The person strata value (H) of 3 suggests that the students are classifiable into three groups of conceptual understanding (high, moderate, and low). On top of that, if the item's separation index value (+9.54) is processed by the same formula (H), the generated value is 13. Such a number shows that the items in the instrument are classifiable into 14 levels of difficulty. Moreover, the data illustrate that the items are deemed accurate and capable of measuring the students' competence in explaining the focused topic.

From the analysis result of students' answer pattern, the research generates Infit and Outfit MNSQ values of 1.02 and 1.05, respectively, with expectation value of 1.0. This clarifies that the students' answer pattern towards the instrument is categorized as "good". In addition, the result generates Infit ZSTD and outfit ZSTD value of 0.0 and 0.10, respectively, with an expectation value of 0.0; the numbers depict that the overall students' answer pattern is in accordance with the model. Moreover, the overall reliability of students section is 0.82, categorized as "good". From the instrument item assessment, it is generated that the Infit and Outfit MNSQ values are 1.03 and 1.05, respectively, with the expectation value of 1.0, and the Infit and Outfit ZSTD values are 0.0 and 0.3, with the expectation value of 0.0. The numbers suggest that the overall instrument is deemed as "good", proven by the instrument reliability value of 0.99. The KR-20 (alpha Cronbach) value results in 0.84, thus signifying a

Commented [A19]: Which table? States the table must be clear.

Commented [A20]: Include the source that states the INFIT and OUTFIT MNSQ categories are Good.

Commented [A21]: What are good categories based on?

Include a reference source or similar research as a reference that indicates the INFIT/OUTFIT ZSTD, and the instrument reliability is categorized as Good.

good interaction between the students and the item. As acquired from the findings, the actual data in this study have met the Rasch model requirements, meaning that further analysis is considered as valid to conduct.

Level of Students' Learning Progress

The second problem of the research is: "How is the learning progress of the participants ranging from senior high school to fourth college year in explaining the focused topic?". To elaborate on that matter, the study employs data generated from the development process of 4TMC instrument to measure the students' conceptual understanding level.



Figure 1 Mean student performance level by grade

(Senior high school students = A, first-year college students = B, second-year college students = C, third-year college students = D, fourth-year college students = E)

Figure 1 displays the average competence calculated in the form of logs based on the students' academic level, ranging from A to E. The figure shows an increasing trend in students' competence development based on their respective academic level (ABCDE). Moreover, it is discovered that the group E shows better learning progress compared to the other groups (D, C, B, and A). Despite that, the One-way ANOVA test indicates a difference among the students' competence based on the academic level, in which F_{count} (6, 0142442) > F_{table} (2,39308); df = 422; p <0.05. The research, therefore, conducted a post hoc Bonferroni test to identify which group that experience significant learning progress. As extracted from the statistical result, group A and B undergo significant learning progress, while group C, D,

and E do not experience such significant advancement. This contradicts the common notion that the group CDE are college students with longer formal education experience compared to group A or B. Such finding indicates that the group CDE find it hard to explain the concept of change of state of matter.

Comparison of average competence between groups ABCDE is conducted to map out the difference in the students' learning progress in each conceptual understanding level (displayed in Table 3). The students' competence is calculated based on four items in each level of conceptual understanding. As an example, in the level 1, the students' competence is measured by referring to the mean of item 1LG-1, 6SL-1, 11SG-1, and 16SL-1; the same also applies in the next levels. Based on Table 4, it is found that the students' competence in level 1 (0.77 logit, SD = 0.86) is higher than their competence in level 2 (0.69 logit, SD = 0.95); the same also applies in the next levels. The findings above indicate that the students' conceptual understanding has not developed optimally. On top of that, the item sequence in level 1 is easier to explain compared to that in level 2. The same condition also applies in the next levels. Students find it harder to explain concepts of change of state of matter as the learning progress level increases. Simply put, the students' learning progress level is different in each level of conceptual understanding.

Table 3

Conceptual	Students' Education Level (Mean, SD)						
Understanding - Level	A (N=171)	B (N=83)	C (N=66)	D (N=55)	E (N=52)	ABDCE (N=427)	
1	0.69 (0.86)	0.80 (0.71)	0.61 (0.91)	1.29 (0.95)	1.05 (0.90)	0.77 (0.86)	
2	0.58 (1.04)	0.66 (0.75)	0.68 (0.86)	1.05 (1.00)	0.83 (0.97)	0.69 (0.95)	
3	0.19 (0.95)	0.61 (1.00)	0.33 (1.13)	0.84 (0.92)	1.10 (1.24)	0.51 (1.10)	
4	0.24 (1.00)	0.53 (0.68)	0.51 (1.12)	0.70 (0.86)	0.51 (0.71)	0.41 (0.57)	
5	-1.16	-0.80	-0.86	-0.48	-0.58	-0.84 (1.41)	

Measurement of students' average competence in each level of conceptual understanding

The difference in students' learning progress levels in each conceptual understanding level depicts that longer formal education experience does not necessarily guarantee that the student will have better learning progress in explaining the focused topic. For instance, Table 4 illustrates the comparison of item logit size in level 3 that is calculated based on the students' academic level.

				•
Table 4 Δx	verage item	logit in	level -	٠
1 abic + 1 i	crage nem	10gn m	IC VCI .	,

Education		Item mean (logit) at level 3				
Level	N	13SG- 3	3LG-3	18LS-3	8SL-3	
А	171	0.51	0.33	-0.22	-0.61	
В	83	0.55	0.40	-0.43	-0.51	
С	66	0.61	0.19	-0.15	-0.66	
D	55	0.33	0.20	-0.15	-0.46	
Е	52	0.57	0.06	-0.30	-0.33	

Discussion

The result shows that: firstly, based on the logit size, the items are put in the following order: 13SG-3 > 3LG-3 > 18 SL-3 > 8SL-3. This is to say that it is harder for the students to explain the concept in item 13SG-3 compared to 3LG-3, 18SL-3, and 8SL-3. Secondly, the students' competence in each item is different and not in sequential order based on the education level (ABCDE). The finding leads to an assumption that all students in group E are supposed to perform better in explaining the item sequence in level 3 than those in group D, C, B, and A, since they progressed through longer education experience. However, the calculation result shows a different insight. In the item 13SG-3, students in group C are the most competent among all group (C (0.61) > E (0.57) > B (0.55) > A (0.51) > D (0.33)), while in the item 8SL-3, group E students are the most competent (E (-0.33) > D (-0.46) > B (-0.51) > A (-0.61) > C (-0.66)). Such a finding indicates that the students' competence is varied despite

being at the same level. To put it another way, longer formal education tends to have an insignificant effect on the development of students' conceptual understanding.

Table 5

Category of item 13SG-3 comprehension

Grade	Ν	Conceptual Understanding Category - Item 13SG-3 (%)				
		LOK	AM	MFN	MFP	SK
А	171	36	21	3	20	19
В	83	19	36	8	5	31
С	66	36	27	2	8	27
D	55	13	24	4	7	53
E	52	23	12	6	12	48

How is the students' learning progress level in the same item? Table 5 displays the percentage data of students' competence in explaining item 13SG-3 based on five categories of conceptual understanding (LOK, AM, MFN, MFP, and SK). In the SK category, students in group D perform better among all groups (D (53%) > E (48%) > B (31%) > C (27%) >A (19%)). Simply put, more than half students in group D are capable of explaining the item 13SG-3 compared to students in other groups. Meanwhile, in LOK, students in group A and C show higher percentage among all groups (A (36%) = C (36%) > E (23%) > B (19%) > D(13%)). In other words, more than one-third of students in group A or C is incapable of explaining the item 13SG-3 compared to students in other groups due to the limited knowledge on the item. Moroever, in AM, group B shows highest percentage among all groups (B (36%) > C (27%) > D (24%) > A (21%) > E (12%)); it signifies that more than one-fourth of students in group B are incapable of explaining item 13SG-3 compared to other groups due to the misconception on the item. Such findings indicate that the high percentage in LOK and AM category is seen as one of the reasons why the students' competence is different in explaining the same item 13SG-3. To put it another way, the students' learning progress does not develop optimally in explaining item 13SG-3 due to lack of knowledge (LOK) or misconception (AM) on the item.



Figure 2(a) Probability Category Curve of item 13SG-3 of group A, and Figure 2(b) Probability Category Curve of item 13SG-3 of group D (Category: LOK = Lack of Knowledge, AM = All-Misconception, MFN = Misconception False Negative, MFP = Misconception False Positive, SK = Scientific Knowledge)

Figure 2 illustrates the comparison of the probability category curve (PCC) of students in group A and D in item 13SG-3. The five curve shapes are the visual representation of the distribution of five categories of students' conceptual understanding. From the curves, one can identify which groups that tend to show LOK and AM category traits. It is worth noting that the curve 2(a) and 2(b) tend to be different based on the MFP curve shape, while others are relatively similar. The MFP curve of students A has a higher probability compared to that of students D; simply put, a senior high school student tends to show stronger MFP category compared to a third-year college student. The notion is supported by the finding that senior high school students are relatively incapable of providing correct reason on item 13SG-3 compared to third-year college students. On the other hand, students with low ability in group D tend to show similar curve shape of LOK, AM, and MFN with group A. This implies that both groups' conceptual understanding in the item is relatively similar. In other words, the learning progress of group D, particularly in students with low ability, has not developed

optimally despite the fact that that group D consists of third-year college students that progressed through three years of formal education experience in university.

This echoes previous findings that the learning progress is highly dependent on the students' learning process and experience (Duschl et al., 2011; Park et al., 2017; Wilson, 2009). Learning progress is defined as a sophisticated and systematic way of thinking, in which the students will undergo gradual progress when learning a topic for a long time interval. Such a systematic way of thinking is formed by the learning practices and education experience (Emden et al., 2018). On top of that, the research findings are in line with previous studies that highlighted that students have distinctive comprehension formed by their own experience (Chi et al., 2018; Emden et al., 2018; Hoe & Subramaniam, 2016; Jin, Mikeska, Hokayem, & Mavronikolas, 2019; Rogat et al., 2011; Testa et al., 2019). Such distinctive knowledge has not been explored by evaluation or intervention through learning roadmaps that are in accordance with remedial learning (Smith et al., 2006). In spite of that, it is considered essential to conduct a further analysis that focuses on the modification of conceptual understanding category and analysis variation that is able to define the characteristics of students' alternative conception. The development procedures, as explained in the methodology, has resulted in 4TMC instrument; however, instrument development is seen as an essential continuous process (Wilson, 2009, 2012).

Based on the research findings, the study identifies several important notes on the development of the 4TMC instrument. Firstly, further analysis of the characteristic of students' response behavior is necessary to conduct regarding the item clarity and the measured concept. The findings have implied that the percentage of LOK and AM understanding category is relatively dominant and tends to increase along with the level of conceptual understanding. Hence, the development of the concept level requires taking into consideration any potential term use that might confuse the students. A further study on the identification of commonly-understood terms or concepts is therefore essential. Secondly, a

separate analysis is required to diagnose the factors contributing to the students' lack of knowledge and misconception. Regarding that, further analysis can be conducted by applying the analysis methods developed by previous studies (Caleon & Subramaniam, 2010; Hoe & Subramaniam, 2016; Oon & Subramaniam, 2013). Thirdly, it is discovered that the concepts LG, SG, SL and LS were interpreted differently by the students. Despite being in the same conceptual understanding level, the items' difficulty level are completely different. Therefore, an evaluation on answer choices requires one to focus on the representation of understanding at the same level.

One of the features of the Rasch model is that the model facilitates one to identify any correlation between the construct map and the students' competence in ways that the students' competence can be analyzed by referring to the difference in item difficulty level. The 4TMC instrument indicates that there are students with very high ability as well as students with low ability in each group. Such a gap serves as the basis for qualitative interpretation to elaborate on the difference in students' competence. The insight is applicable in the learning process of chemistry subject. The instrument is expected to be beneficial for teachers in developing a formative test to identify the students' progress of conceptual understanding. On top of that, teachers are able to implement the instrument as a diagnostic instrument to evaluate students' conceptual understanding in providing feedback on their learning progress. Further, the teachers will be able to develop instructional strategies that are specifically designed to tackle the students' difficulty in developing an epistemological explanation regarding the concept of change of state of matter.

Commented [A22]: Why? You need to compare the findings with relevant theories.

Conclusions

The result revealed that the integration of the 4TMC test and Rasch modeling is effective and valid to be treated as the diagnostic instrument to measure students' learning progress. Moreover, it is discovered that students in group A, B, C, D, and E, particularly those with

Commented [A23]: Conclusions must be able to answer the research purposes and be short enough.

low ability, are hampered in developing an epistemological explanation of the concept. This blames the students' lack of certainty in their answer and reason; thus, assumed as having lack of knowledge or misconception. The low-ability students' curve shape of LOK and AM is consistent in the competence interval of less than 0.1 logit. On the other hand, the students' ability gets lower as the conceptual understanding level increases. Such finding indicates that the learning process and education experience provide a limited contribution for the students in developing a systematic way of thinking regarding the concept of change of state of matter.

Recommendations

The Based on the results of the study, there are several recommendations for researchers and teachers. For researchers, the findings of this research can be followed up to examine more in how students build their understanding gradually in explaining the concept of particles in substance form changes. The study can be conducted by developing tests that aim to evaluate and diagnose the process of student knowledge formation and development while being able to identify at the level of education where the confusion of understanding occurs. The evaluation becomes more objective, not only reviewed from the student's point of ability but can be reviewed from the teacher's ability. The model of *Rasch's* multi-faced item response pattern approach becomes one of the important parts recommended for such objectives. In this way, students' ability to develop epistemological knowledge, and their ability to significantly actualize the knowledge gained can be measured well.

On the other hand, for teachers, the results of this study along with the stages of analysis approach used can be a reference in evaluating the progress of learners' learning, as well as determining alternative thinking frameworks of students in explaining the concept of substance change. The information serves as strategic feedback in formulating instructional strategies and preparing remedial learning, especially for students who have difficulty in developing epistemological explanations of substance changes.

Limitations

The limitations of the research are primarily related to the misrepresentation of student reasoning, which may arise in its efforts to connect phenomena and concepts measured in each item. In this context, the student may not excel to explain, because of his incapableness in using his heuristic reasoning. This instrument is not equipped with items that evaluate the heuristic abilities of the student in question. However, researchers decided to record this incompetence as a misconception or vague knowledge. For further research, it is recommended that the instrument be equipped with items that measure students' emotional and heuristic reasoning according to the conceptual framework to be evaluated.

Acknowledgments

The researchers would like to express their gratitude towards the Directorate of Research and Community Service, Ministry of Research and Technology of Republic of Indonesia, for the financial support through the University Basic Research Excellence Grant Program in the Research and Community Service Office of Universitas Negeri Gorontalo, 2020.

References

- [Aktan, D. C. (2013). Investigation of students' intermediate conceptual understanding levels: The case of direct current electricity concepts. *European Journal of Physics*, 34(1), 33– 43. https://doi.org/10.1088/0143-0807/34/1/33
- Arslan, H. O., Cigdemoglu, C., & Moseley, C. (2012). A three-tier diagnostic test to assess pre-service teachers' misconceptions about global warming, greenhouse effect, ozone layer depletion, and acid rain. *International Journal of Science Education*, 34(11), 1667– 1686. https://doi.org/10.1080/09500693.2012.680618
- Bond, T. G., & Fox, C. M. (2007). *Applying The Rasch Model: Fundamental Measurent in the Human Sciences* (2nd Ed.). Routledge Taylor & Francis Group

- Caleon, I. S., & Subramaniam, R. (2010). Do students know what they know and what they don't know? Using a four-tier diagnostic test to assess the nature of students' alternative conceptions. *Research in Science Education*, 40(3), 313–337. https://doi.org/10.1007/s11165-009-9122-4
- Chandrasegaran, A. L., Treagust, D. F., & Mocerino, M. (2007). The development of two tier multiple-choice diagnostic instrument for evaluating secondary school students' ability to describe and explain chemical reactions using multiple levels of representation. *Chemistry Education Research and Practice*, 8(3), 293–307
- Chi, S., Wang, Z., Luo, M., Yang, Y., & Huang, M. (2018). Student progression on chemical symbol representation abilities at different grade levels (Grades 10–12) across gender. *Chemistry Education Research and Practice*, 19(4), 1055–1064. https://doi.org/10.1039/c8rp00010g
- Claesgens, J., Scalise, K., Wilson, M., & Stacy, A. (2009). Mapping student understanding in chemistry: The perspectives of chemists. *Science Education*, 93(1), 56–85. https://doi.org/10.1002/sce.20292.
- Duncan, R. G., & Hmelo-Silver, C. E. (2009). Learning progressions: Aligning curriculum, instruction, and assessment. *Journal of Research in Science Teaching*, 46(6), 606–609. https://doi.org/10.1002/tea.20316
- Duschl, R., Maeng, S., & Sezen, A. (2011). Learning progressions and teaching sequences: A review and analysis. *Studies in Science Education*, 47(2), 123–182. https://doi.org/10.1080/03057267.2011.604476
- Emden, M., Weber, K., & Sumfleth, E. (2018). Evaluating a learning progression on "Transformation of Matter" on the lower secondary level. *Chemistry Education Research and Practice*, 19(4), 1096–1116. https://doi.org/10.1039/c8rp00137e
- Habiddin, & Page, E. M. (2019). Development and validation of a four-tier diagnostic instrument for chemical kinetics (FTDICK). *Indonesian Journal of Chemistry*, 19(3), 720–736. https://doi.org/10.22146/ijc.39218
- Hadenfeldt, J. C., Bernholt, S., Liu, X., Neumann, K., & Parchmann, I. (2013). Using ordered multiple-choice items to assess students' understanding of the structure and composition of matter. *Journal of Chemical Education*, 90(12), 1602–1608. https://doi.org/10.1021/ed3006192
- Hasan, S., Bagayoko, D., & Kelley, E. L. (1999). Misconceptions and the certainty of response index (CRI). *Physics Education*, 34(5), 294–299. https://doi.org/10.1088/0031-9120/34/5/304
- Herrmann-Abell, C. F., & Deboer, G. E. (2016). Using rasch modeling and option probability curves to diagnose students' misconceptions. *American Eduacational Research Assossiation*, 1–12
- Hoe, K. Y., & Subramaniam, R. (2016). On the prevalence of alternative conceptions on acidbase chemistry among secondary students: Insights from cognitive and confidence measures. *Chemistry Education Research and Practice*, 17(2), 263–282.

https://doi.org/10.1039/c5rp00146c

- Jin, H., Mikeska, J. N., Hokayem, H., & Mavronikolas, E. (2019). Toward coherence in curriculum, instruction, and assessment: A review of learning progression literature. *Science Education*, 103(5), 1206–1234. https://doi.org/10.1002/sce.21525
- Klassen, S. (2006). Contextual assessment in science education: Background, issues, and policy. *Science Education*, 90(5), 820–851. https://doi.org/10.1002/sce.20150
- Laliyo, Botutihe, & Panigoro. (2019). The development of two-tier instrument based on distractor to assess conceptual understanding level and student misconceptions in explaining redox reactions. *International Journal of Learning, Teaching and Educational Research*, 18(9), 216–237. https://doi.org/10.26803/ijlter.18.9.12
- Linacre, J. M. (2012). A user's guide to WINSTEPS ® MINISTEP Rasch-model computer program: Program manual 3.75.0. https://doi.org/ISBN 0-941938-03-4
- Linacre, J. M. (2020). A User's Guide to WINSTEPS ® MINISTEP Rasch-Model Computer Programs Program Manual 4.5.1. https://doi.org/ISBN 0-941938-03-4
- Ling Lee, W., Chinna, K., & Sumintono, B. (2020). Psychometrics assessment of HeartQoL questionnaire: A Rasch analysis. *European Journal of Preventive Cardiology*, 1–6. https://doi.org/10.1177/2047487320902322
- Liu, X. (2012). Developing measurement instruments for science education research. In B. Fraser, K. G. Tobin, & C. J. McRobbie (Eds.), *Second international handbook of science education* (pp. 651–665). Springer Netherlands
- Lu, S., & Bi, H. (2016). Development of a measurement instrument to assess students' electrolyte conceptual understanding. *Chemistry Education Research and Practice*, 17(4), 1030–1040. https://doi.org/10.1039/c6rp00137h
- Morell, L., Collier, T., Black, P., & Wilson, M. (2017). A construct-modeling approach to develop a learning progression of how students understand the structure of matter. *Journal of Research in Science Teaching*, 54(8), 1024–1048. https://doi.org/10.1002/tea.21397
- Neumann, K., Viering, T., Boone, W. J., & Fischer, H. E. (2013). Towards a learning progression of energy. *Journal of Research in Science Teaching*, 50(2), 162–188. https://doi.org/10.1002/tea.21061
- Park, M., Liu, X., & Waight, N. (2017). Development of the connected chemistry as formative assessment pedagogy for high school chemistry teaching. *Journal of Chemical Education*, 94(3), 273–281. https://doi.org/10.1021/acs.jchemed.6b00299
- Peterson, R. F., Treagust, D. F., & Garnett, P. (1989). Development and application of a diagnostic instrument to evaluate grade-11 and -12 students' concepts of covalent bonding and structure following a course of instruction. *Journal of Research in Science Teaching*, 26(4), 301–314. https://doi.org/10.1002/tea.3660260404
- Rogat, A., Anderson, C., Foster, J., Goldberg, F., Hicks, J., Kanter, D., ... Wiser, M. (2011). Developing learning progression in support of the new science standards: A RAPID

workshop series. (4), 163. https://doi.org/10.12698/cpre.2011.lprapid

- Sadler, P. M. (1999). The relevance of multiple-choice testing in assessing science understanding. In J. J. Mintzes, J. H. Wandersee, & J. D. Novak (Eds.), Assessing science understanding: A human constructivist view (pp. 251–274). Zaccheus Onumba Dibiaezue Memorial Libraries. https://zodml.org/sites/default/files/%5BJoel_J._Mintzes%2C_James_H._Wandersee%2 C_Joseph_D._No_0.pdf
- Smith, C. L., Wiser, M., Anderson, C. W., & Krajcik, J. (2006). Implications of research on children's learning for standards and assessment: A proposed learning progression for matter and the atomic-molecular theory. *Measurement: Interdisciplinary Research & Perspective*, 4(1–2), 1–98. https://doi.org/10.1080/15366367.2006.9678570
- Sumintono, B., & Widhiarso, W. (2014). Aplikasi model Rasch untuk penelitian ilmu-ilmu sosial [Application of Rasch model in social science research]. Trim Komunikata. https://www.researchgate.net/publication/268688933%0AAplikasi
- Testa, I., Capasso, G., Colantonio, A., Galano, S., Marzoli, I., Scotti di Uccio, U., ... Zappia, A. (2019). Development and validation of a university students' progression in learning quantum mechanics through exploratory factor analysis and Rasch analysis. *International Journal of Science Education*, 41(3), 388–417. https://doi.org/10.1080/09500693.2018.1556414
- Treagust, D. F. (1988). Development and use of diagnostic tests to evaluate students' misconceptions in science. *International Journal of Science Education*, 10(2), 159–169. https://doi.org/10.1080/0950069880100204
- Tyson, L., Treagust, D. F., & Bucat, R. B. (1999). The complexity of teaching and learning chemical equilibrium. *Journal of Chemical Education*, 76(2–4), 554–558. https://doi.org/10.1021/ed077p1560.1
- Wilson, M. (2005). Constructing measures: an item response modeling approach. Lawrence Erlbaum Associates, Inc. https://doi.org/10.4324/9781410611697
- Wilson, M. (2008). Cognitive diagnosis using item response models. Zeitschrift Für Psychologie / Journal of Psychology, 216(2), 74–88. https://doi.org/10.1027/0044-3409.216.2.74
- Wilson, M. (2009). Measuring progressions: Assessment structures underlying a learning progression. *Journal of Research in Science Teaching*, 46(6), 716–730. https://doi.org/10.1002/tea.20318
- Wilson, M. (2012). Responding to a challenge that learning progressions pose to measurement practice. In A. C. Alonzo & A. W. Gotwals (Eds.), *Learning progression in science* (pp. 317–344). https://doi.org/10.1007/978-94-6091-824-7

APPENDIX

Commented [A24]: 1.Citation on behalf of Linacre, J.M (2020) is not listed. Please check again.

2. Sadler (1999) has a year that does not match the contents of the paper.

3.Citation on behalf of Sumintono & Widhiarso have different years in the content of the paper, ie (2014) and (2015), but the references are not listed. Please check again.

4. What is the difference between Linacre (2012) and Linacre (2020) apart from the edition and program manual version?

5.The majority of references are recent. Few are old (1988, 1989, 1999). There are few references in this paper. Add the latest supporting references indexed by Scopus or reputable journals.

Commented [A25]: Add Appendix and insert 4TMC instrument.

Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress

Abstract: The purpose of this study was to developing and implementation four-tier multiple-choice (hereinafter, 4TMC) instrument with Rasch model to evaluate students' learning progress in explaining the concept of change of state of matter. The data were obtained through development and validation techniques on 20 4TMC items distributed to 427 students. On each item, the study applied diagnostic-summative assessment and certainty response index. The students' conceptual understanding level was categorized based on the combination their answer choices; the measurement generated Partial-Credit polytomous Rasch model data. The data were further processed by WINSTEPS version 4.5.3 software to equate the data interval rate. Analysis of differences based on class level of students using Analysis of Variants (One-way ANOVA). The result revealed that the integration of 4TMC test and Rasch modeling was effective to be treated as the instrument to measure students' learning progress. One-way ANOVA test indicates a difference among the students' competence based on the academic level, in which F_{count} (6, 0142442) > F_{table} (2,39308); df = 422; p <0.05. On top of that, it was discovered that low-ability students see very slow progress due to the lack of knowledge as well as a misconception in explaining the concept as mentioned above.

Keywords: Learning progress. four-tier, change of state of matter, Rasch model.

Introduction

Central to the notion of science learning is the development of students' scientific understanding of basic concepts of sciences (Hadenfeldt et al., 2013), particularly, change of state of matter (Emden et al., 2018). Aside from the issue, several studies have also highlighted the students' inability to provide an epistemological explanation of basic concepts of sciences (Chi et al., 2018). Efforts to solve the issues, however, have shown little progress, as the students might have more complex perceptions regarding the alternative concept they understand (Morell et al., 2017).

Education practitioners have recommended the utilization of learning progress concept as the instructional method to provide guidance and direction and to adjust the curriculum with the learning process and assessment (Claesgens *et al.*, 2009; Duncan & Hmelo-Silver, 2009;

Rogat *et al.*, 2011). Learning progress is defined as a sophisticated and systematic way of thinking. This method is applicable for a learning process, in which the students will undergo gradual progress when learning a topic in a long duration. Its effectiveness is highly dependent on the learning process and the students' learning experience (Duschl et al., 2011). The concept involves certain sets of gradual levels that represent conceptual understanding, ranging from low level up to comprehensive level.

The notion of learning progress is highly distinctive to each student and is dependent to one's learning experience (Rogat et al., 2011); therefore, there is no learning roadmap that is suitable for all kinds of students (Smith et al., 2006). Each student constructs one's understanding in a different way; moreover, the construction process is varied depending on the students' conceptual understanding level (Aktan, 2013). This is to say that each student undergoes a different rate of learning progress, understanding level, and knowledge construction. Simply put, the development of scientific comprehension among students is not linear (Neumann et al., 2013). Therefore, this study regards each level of students' conceptual understanding for more advanced level of understanding (Hadenfeldt et al., 2013). A student who faces difficulty in a certain level of understanding will see a lack of progress to a more advanced level. This in turn hinders the student's ability to construct an epistemological explanation on the basic concepts of science.

Efforts to diagnose the epistemological problems, as mentioned previously, are feasible to conduct if the extent of students' conceptual understanding is formulated. Within this context, the learning progress is treated as the method to evaluate students' conceptual understanding. The diagnostic information generated is reliable to be treated as a reference for the teachers in developing accurate and valid instructional components to guide the students to progress to the next level. Despite the potentials, this study deems that it is challenging for the teachers to construct such an accurate instrument.

Among the diagnostic instruments that are considered applicable is the 4TMC instrument. It is the development of two-tier multiple-choice test recommended by Treagust (1988) and Chandrasegaran et al., (2007). The use of two-tier instrument is familiar in identifying students' understanding in select topics such as electrochemistry (Lu & Bi, 2016), covalent bond (Peterson, Treagust, & Garnett, 1989), and chemical equilibrium (Tyson et al., 1999). Despite its reputation in academia, the two-tier test has raised criticism due to its sole focus on the facts and negligence towards students' understanding (Klassen, 2006). Therefore, several experts propose the renewed version of the test by adding distractor answer choices to strengthen the diagnostic value of the items (Herrmann-Abell & DeBoer, 2011; Sadler, 1998). In addition, some have highlighted the test's weakness in cases where students' tended to pick the answer choice and the reasoning randomly. This illustrates that the students were uncertain and possessed several misconceptions in the first tier question. In such cases, teachers faced difficulty in differentiating between guessed answers and misconceptions (Habiddin & Page, 2019; Hasan et al., 1999).

The criticism laid against the model has sparked the innovation of three-tier and four-tiers instruments. Both instruments feature two multi-level questions, also similar with two-tier test. In the three-tier test, however, the measurement of students' certainty level is conducted simultaneously in both first and second-tier questions; in the meantime, the measurement is conducted separately in the first two tiers (Caleon & Subramaniam, 2010). The value of students' certainty rate ranges from one (very uncertain) to five (very certain).

Three-tier test lacks validity in measuring the students' certainty rate regarding both the answer choice and the reasoning, whether or not the value of certainty rate refers only to the answer choice, to the reasoning, or both. Such weakness will in turn obstructs the evaluation and classification process of students' responses (Arslan et al., 2012). In the four-tier instrument, the measurement of certainty rate also involves the answer choice in the first tier and the reasoning in the third tier (Arslan et al., 2012; Loh et al., 2014). Regarding this

Commented [MOU1]: Revise it according to APA 7.

feature, four-tier test is considered more accurate than the three-tier test. Students who pick wrong answer choices with high certainty indicate that they have a very high misconception on the measured item (Hoe & Subramaniam, 2016).

Four-tier instruments are used in studies discussing topics such as physics education (Caleon & Subramaniam, 2010), chemical thermodynamics (Sreenivasulu & Subramaniam, 2013), transition metal (Sreenivasulu & Subramaniam, 2014), acid-base reaction (Hoe & Subramaniam, 2016), and chemical kinetics (Habiddin & Page, 2019). However, it is worth noticing that studies on chemistry topic which employ four-tiers instruments tend to focus on describing alternative conception. To put it another way, the higher the certainty rate is, the stronger the students' alternative conception will be. Despite its potentials, the scholarly discussion has overlooked the implementation of a four-tier diagnostic instrument to measure students' learning progress. Therefore, further analysis is essential on the application of 4TMC test in several domains analyzes by Rasch model approach.

The use of Rasch model has been introduced since the 2000s in the science education research; it features the instrument that integrates diagnostic assessment and summative assessment (Liu, 2012; Wei et al., 2012). On top of that, the diagnostic assessment approach is introduced to conduct an in-depth analysis of the construction process of students' conceptual understanding (Claesgens et al., 2009; Hadenfeldt et al., 2013; Lu & Bi, 2016). This study employs 4MTC test and Rasch modeling as a diagnostic tool to evaluate students' learning progress in explaining the change of state of matter. The study focus revolves around two research questions: 1) How is the effectiveness of 4TMC instrument to evaluate the students' learning progress in explaining concepts of change of state of matter. 2) How is the learning progress in students ranging from the senior high school level up to the senior (fourth) year of college in explaining the concepts?

Methodology

Development Model

This research used a development research referring to the test development model from Wilson. Wilson (2005, 2008) introduces four steps of measurement instrument development: The first step is to the learning progress variable focused on a characteristic measured at a particular time unit. The second step comprises the design process of items or tasks used to measure students' responses. Moreover, the third step involves outcome space, in which the students' responses are categorized into all items related with the learning progress variable. On top of that, the fourth step employs measurement model, such as Rasch model. This recommendation is proven valid to be implemented in developing measurement instrument for different construct variables (Barbera, 2013; Chi et al., 2018; Hadenfeldt et al., 2013; Laliyo, Botutihe, & Panigoro, 2019; Lu & Bi, 2016; Wei et al., 2012; Wilson, 2009; Wind, Tsai, Grajeda, & Bergin, 2018). The study conducted development of measurement instruments by referring to Wilson's recommendation (2005, 2008) and adopted Treagust's framework (1988) of item development. The present study also included two questions related to certainty rate (Arslan et al., 2012; Habiddin & Page, 2019: Hasan et al., 1999). The obtained data were analyzed by Rasch model approach.

Construct Map: Determining Level of Understanding

The first step was to develop the construct of measured variables. The study involved four concepts of change of state of matter: liquid-gas (LG), solid-liquid (SL), solid-gas (SG), and liquid-solid (LS). These concepts were implemented in a gradual manner through five levels of conceptual understanding (Table 1). Such method functions as the pathway of conceptual development that involves learning objectives from the lowest to the highest level of conceptual understanding (Duncan & Hmelo-Silver, 2009; Löfgren & Helldén, 2009; Hadenfeldt et al., 2013; Rogat et al., 2011). In other words, the set of levels, as mentioned previously, was adjusted to the students' needs so as to develop their conceptual

Commented [MOU2]: Why did you ignore gas-solid(GS) and gas-liquid(GL)?

understanding. This took into account that each student might progress on different and nonlinear development of conceptual understanding; therefore, the levels, as illustrated in Table 1, was considered valid to illustrate the ideal conceptual development pathway (Neumann et cl. 2012).

al., 2013).

Table 1. Level of Conceptual Understanding in Explaining Concept of Change of State of Matter

Conceptual Understanding Level		Change of State of Matter/Item				
		LG	SL	SG	LS	
5	5 Submicroscopic diagram of change of		10SL-5	15SG-5	20LS-5	
	state of matter					
4	Correlation between state of matter and	4LG-4	9SL-4	14SG-4	19LS-4	
	the process of change of state of matter					
3	Process of change of state of matter	3LG-3	8SL-3	13SG-3	18LS-3	
2	Concept of state of matter	2LG-2	7SL-2	12SG-2	17LS-2	
1	Factual phenomenon of state of matter	1LG-1	6SL-1	11SG-1	16LS-1	
examination: (I.C liquid gas, SI solid liquid SC solid gas, I.S liquid gas)						

Description: (LG = liquid-gas, SL = solid-liquid, SG = solid-gas, LS = liquid-gas)

Item Design and Assessment Scheme

The second phase involved an item design. In the 4TMC instrument, all the items consisted of four-tier multiple-choices. To put it another way, each item contains four questions that combine between diagnostic-summative test (Hoe & Subramaniam, 2016; Lu & Bi, 2016; Treagust, 1988) with certainty response index (hereinafter, CRI) test (Arslan et al., 2012; Hasan et al., 1999). The first-tier questions (Q1) aimed to identify whether or not the students understand the content. Moreover, questions in the second tier (Q2) were employed to clarify the students' certainty regarding their answers in the Q1. Third-tier questions (Q3) functioned to diagnose the students' reasoning regarding their answers in the Q1. Further, questions in the second tier (Q4) were employed to clarify the students' certainty regarding their answers in the Q3. Q1 and Q3 questions in each item involved five answer choices; one among them was the correct answer, while three were the distractor, and another answer choice was open-ended answer choice. This open-ended option allows the students to decide the answer by themselves, should they find no correct answer as in accordance with their conceptual understanding. In the meantime, the Q2 and Q4 questions involved two close-ended answer choices; the first choice was for those who are uncertain of their answer, and the second
choice was for the students who are very certain of their answer (Arslan et al., 2012). The distractor choices were employed in Q1 and Q3 questions to validate the diagnostic strength of the questions (Herrmann-Abell & DeBoer, 2011; Sadler, 1998). Therefore, in the Q1 and Q3 tiers, the students would have only 0.20 or 20 percent probability of choosing the correct answer.

Outcome Space and Data Collection

The third step involved the design of the outcome space of the correlation between items and construct maps (Bond & Fox, 2007; Wilson, 2009). The item validation was conducted independently by three expert validators to evaluate the extent of correlation between answer choices in Q1-Q3 in each item and the level of students' conceptual understanding. The validators were asked to clarify that the questions are easy to understand and the students' lack of linguistic competence would not hinder them from providing the right answer. The validators also required to ensure that the questions are in accordance with the syllabus, particularly with the students' conceptual understanding as based on the construct map. The questions in each item were also validated in several aspects, such as: ambiguity, time allocation, directiveness towards a particular answer, and subjective or emotional expression. Fleiss κ measure was employed to acquire information on the validators' approval. From the measure, it was generated that the κ value = 0.97, indicating that the three validators agreed that the 4TMC items were valid in correlating between the answer choices and the students' conceptual understanding.

The next step was to acquire data based on the measurement instrument. The instrument was distributed to 427 students in Gorontalo, Indonesia. The students comprised 171 (40.05%) senior high school students (or students A), 83 (19.44%) university freshmen majoring chemistry education (or students B), 66 (15.45%) second-year university students majoring chemistry education (or students C), 55 (12.88%) third-year university student majoring chemistry education (or students D), and 52 (12.18%) fourth-year university students

majoring chemistry education (or students E). Based on gender, the female participants comprised 369 participants (86.41%), and the male counterparts consisted of 58 participants (13,58%). The participants were given no particular educational treatments and had stated their voluntary consent to participate in the research.

Rasch Model Measurement and Data Analysis

The fourth step was to conduct the Rasch model measurement. This step was implemented to define the correlation between the score generated and the students' conceptual understanding level as elaborated within the construct map. The involvement of Rasch model measurement lay on the assumption that the item difficulty level is dependent on the students' answer, and that the students' understanding is dependent on the estimation of item difficulty (Linacre, 2012).

Rasch partial credit model (PCM) was employed to evaluate the learning progress through structured questions; this took into account that the instrument items involved gradual and structured questions (Bond and Fox, 2007; Masters, 1982; Sumintono and Widhiarso, 2015; 2009). The Wilson, model was stated into the following formula: $ln[P_nik/(1-P_nik)] - B_n - D_ik$, in which P_{nik} refers to the probability of student n with Bn ability to pick correct response in the level k of item i; while Dik refers to the difficulty level k of item i, or the threshold point for the test taker who scores k, not k -1. Analysis of differences based on class level of students using One-way ANOVA.

Results and Discussion

Results

Effectiveness of Measurement Instruments

Unidimensionality is an essential indicator to evaluate the 4TMC instrument's ability to measure students' capability of explaining the concept of change of state of matter. This indicator is measured by Principal Component Analysis of the residuals to estimate the extent

of variance to which the instrument is able to measure what it is supposed to measure (Sumintono & Widhiarso, 2014). he result of raw variance explained by measures of data is 38.9%, the number almost approaches the expectation value of 39.2%. The numbers indicate that the minimum unidimensionality requirements of 20% are achieved, and simultaneously, the limit of Rasch unidimension is met (approaching 40%) (Linacre, 2012; Ling Lee, Chinna, & Sumintono, 2020). Moreover, the instrument's unexplained variance values are below 7% and considered as ideal (not exceeding 15%), signifying that the item independence rate in instrument falls into "good" category.

The second step is to measure the consistency between the item difficulty level and students' conceptual understanding. The research discovers several interesting cases regarding the difference between the items and students' conceptual understanding: Firstly, there are four items identified (LG, SL, SG and LS) that measure similar constructs within each level of conceptual understanding. Despite being in the same conceptual understanding level, the items' logit is completely different. For instance, four items were discovered in level 3, each with varying logit (8SL-3 (-0.33) < 18LS-3 (-0.29) < 3LG-3 (+0.15) < 13SG-3 (+0.30)). The numbers indicate that overall, students are more capable of explaining the concept of SL state change compared to LS, LG, and SG. This condition also occurs in the level 4, in which each item has varying logit (19LS-4 (-0.37) < 9SL-4 (+0.04) = 14SG-4(+0.04) < 4LG-4 (+0.07)). Such a finding shows that the students find it easier to explain the correlation between the state of matter and the change process in LS compared to either SL, SG, or LG. Two sample cases above have illustrated that the students' conceptual understanding differs between the change process of LG (evaporation), SG (sublimation), SL (melting), and LS (freezing).

Moreover, it is found that the items in higher conceptual understanding levels tend to have lower logit than those at a lower level. As an instance, the logit of item 19SL-4 in level 4 (-0.37) is smaller than that of item 13SG-3 in level 3 (+0.30). This signifies that students find it harder to explain the item 13SG-3 compared to item 19SL-4. Thirdly, in the same concept of **Commented [MOU3]:** Revise it according to APA 7. (Ling Lee at al., 2020)

change of state (for example, LS), the logit of item 17LS-2 in level 2 (-0.40) is smaller than that of item 16LS-1 in level 1 (-0.16). As illustrated by the number, students find it easier to explain the SL concept in level 2 rather than to explain the concept's macroscopic fact in level 1. The findings above indicate that the students' conceptual understanding is not consistent with the item sequence. Moreover, the findings also suggest that the item difficulty level (LG, SL, SG, and LS; particularly SL (melting) and LS (freezing)) do not match the level in the construct map.

Measurement reliability

In Rasch analysis, the indicator of reliability is observed from the quality of students' response patterns, the instrument, and the interaction between person-item. Within this study, item separation and person separation values are employed as the indicators. The separation index is also converted to Cronbach-equivalent value with an estimation of 0-1. The summary of measurement instrument statistics is displayed in Table 2 as follows:

Table 2. Summary of fit statistics

	Student	Item
	(N=427)	(N=20)
Mean	0.26	0.00
Standard Error	0.02	0.09
Standard Deviation (SD)	0.48	0.41
Reliability	0.82	0.99
Infit mean-square	1.02	1.03
Outfit mean-square	1.05	1.05
Infit ZSTD	0.00	0.00
Outfit ZSTD	0.10	0.30
Point Raw Score to measure correlation	0.99	-0.99
Separation index (reliability)	2.10	9.54
Cronbach Alpha (KR-20): 0.84		
Data Points : 8540		
Chi-Square : 21173		
df : 8091 (p = 0.0000)		

From the table 2, it is generated that the total data points are 8540 with a Chi-square value of 21173 and the degree of freedom (df) of 8091 (p = 0.0000). These numbers indicate that the measurement is deemed as "very good" and "significant". The column of students and item in

the table suggest whether or not the students and the item are considered fit. The average measure value of students is +0.26 logit ($\mu > 0.00$), signifying that the students in overall are competent to explain the concept of change of state of matter. If the separation index value of students (+2.10 logit) is inputted into the person strata (H) formula, or H = [(4*separation) + 1]/3, thus, the generated H value = +3.13 (Linacre, 2012; Sumintono & Widhiarso, 2015). The person strata value (H) of 3 suggests that the students are classifiable into three groups of conceptual understanding (high, moderate, and low). On top of that, if the item's separation index value (+9.54) is processed by the same formula (H), the generated value is 13. Such a number shows that the items in the instrument are classifiable into 14 levels of difficulty. Moreover, the data illustrate that the items are deemed accurate and capable of measuring the students' competence in explaining the focused topic.

From the analysis result of students' answer pattern, the research generates Infit and Outfit MNSQ values of 1.02 and 1.05, respectively, with expectation value of 1.0. This clarifies that the students' answer pattern towards the instrument is categorized as "good". In addition, the result generates Infit ZSTD and outfit ZSTD value of 0.0 and 0.10, respectively, with an expectation value of 0.0; the numbers depict that the overall students' answer pattern is in accordance with the model. Moreover, the overall reliability of students section is 0.82, categorized as "good". From the instrument item assessment, it is generated that the Infit and Outfit MNSQ values are 1.03 and 1.05, respectively, with the expectation value of 1.0, and the Infit and Outfit ZSTD values are 0.0 and 0.3, with the expectation value of 0.0. The numbers suggest that the overall instrument is deemed as "good", proven by the instrument reliability value of 0.99. The KR-20 (alpha Cronbach) value results in 0.84, thus signifying a good interaction between the students and the item. As acquired from the findings, the actual data in this study have met the Rasch model requirements, meaning that further analysis is considered as valid to conduct.

Level of Students' Learning Progress

The second problem of the research is: "How is the learning progress of the participants ranging from senior high school to fourth college year in explaining the focused topic?". To elaborate on that matter, the study employs data generated from the development process of 4TMC instrument to measure the students' conceptual understanding level.



Figure 1 Mean student performance level by grade

(Senior high school students = A, first-year college students = B, second-year college students = C, third-year college students = D, fourth-year college students = E)

Figure 1 displays the average competence calculated in the form of logs based on the students' academic level, ranging from A to E. The figure shows an increasing trend in students' competence development based on their respective academic level (ABCDE). Moreover, it is discovered that the group E shows better learning progress compared to the other groups (D, C, B, and A). Despite that, the One-way ANOVA test indicates a difference among the students' competence based on the academic level, in which F_{count} (6, 0142442) > F_{table} (2,39308); df = 422; p <0.05. The research, therefore, conducted a post hoc Bonferroni test to identify which group that experience significant learning progress. As extracted from the statistical result, group A and B undergo significant learning progress, while group C, D, and E do not experience such significant advancement. This contradicts the common notion that the group CDE are college students with longer formal education experience compared to group A or B. Such finding indicates that the group CDE find it hard to explain the concept of change of state of matter.

Commented [MOU4]: Did you calculate from Table 3? They are different from the average of 5 levels.

Comparison of average competence between groups ABCDE is conducted to map out the difference in the students' learning progress in each conceptual understanding level (displayed in Table 3). The students' competence is calculated based on four items in each level of conceptual understanding. As an example, in the level 1, the students' competence is measured by referring to the mean of item 1LG-1, 6SL-1, 11SG-1, and 16SL-1; the same also applies in the next levels. Based on Table 4, it is found that the students' competence in level 1 (0.77 logit, SD = 0.86) is higher than their competence in level 2 (0.69 logit, SD = 0.95); the same also applies in the next levels. The findings above indicate that the students' conceptual understanding has not developed optimally. On top of that, the item sequence in level 1 is easier to explain compared to that in level 2. The same condition also applies in the next levels. Students find it harder to explain concepts of change of state of matter as the learning progress level increases. Simply put, the students' learning progress level is different in each level of conceptual understanding.

Table 3

Conceptual Understanding	Students' Education Level (Mean, SD)					
Level	A (N=171)	B (N=83)	C (N=66)	D (N=55)	E (N=52)	ABDCE (N=427)
1	0.69 (0.86)	0.80 (0.71)	0.61 (0.91)	1.29 (0.95)	1.05 (0.90)	0.77 (0.86)
2	0.58 (1.04)	0.66 (0.75)	0.68 (0.86)	1.05 (1.00)	0.83 (0.97)	0.69 (0.95)
3	0.19 (0.95)	0.61 (1.00)	0.33 (1.13)	0.84 (0.92)	1.10 (1.24)	0.51 (1.10)
4	0.24 (1.00)	0.53 (0.68)	0.51 (1.12)	0.70 (0.86)	0.51 (0.71)	0.41 (0.57)
5	-1.16(1.59)	-0.80(1.46)	-0.86(1.51)	-0.48(0.85)	-0.58(1.51)	-0.84 (1.41)

Measurement of students' average competence in each level of conceptual understanding

The difference in students' learning progress levels in each conceptual understanding level depicts that longer formal education experience does not necessarily guarantee that the student will have better learning progress in explaining the focused topic. For instance, Table 4 illustrates the comparison of item logit size in level 3 that is calculated based on the students' academic level.

Table 4 Average item logit in level 3

Education	N	Iteı	n mean (lo	git) at leve	el 3
Level	19	13SG-3	3LG-3	18LS-3	8SL-3
А	171	0.51	0.33	-0.22	-0.61
В	83	0.55	0.40	-0.43	-0.51
С	66	0.61	0.19	-0.15	-0.66
D	55	0.33	0.20	-0.15	-0.46
E	52	0.57	0.06	-0.30	-0.33

Discussion

The result shows that: firstly, based on the logit size, the items are put in the following order: 13SG-3 > 3LG-3 > 18 SL-3 > 8SL-3. This is to say that it is harder for the students to explain the concept in item 13SG-3 compared to 3LG-3, 18SL-3, and 8SL-3. Secondly, the students' competence in each item is different and not in sequential order based on the education level (ABCDE). The finding leads to an assumption that all students in group E are supposed to perform better in explaining the item sequence in level 3 than those in group D, C, B, and A, since they progressed through longer education experience. However, the calculation result shows a different insight. In the item 13SG-3, students in group C are the most competent among all group (C (0.61) > E (0.57) > B (0.55) > A (0.51) > D (0.33)), while in the item 8SL-3, group E students are the most competent (E (-0.33) > D (-0.46) > B (-0.51) > A (-0.61) > C (-0.66)). Such a finding indicates that the students' competence is varied despite being at the same level. To put it another way, longer formal education tends to have an insignificant effect on the development of students' conceptual understanding.

Table 5

Category of item 13SG-3 comprehension

	_	Conceptu	al Unders	tanding Ca	tegory - Iter	n 13SG-3
Grade	N			(%)		
		LOK	AM	MFN	MFP	SK
A	171	36	21	3	20	19
В	83	19	36	8	5	31
С	66	36	27	2	8	27
D	55	13	24	4	7	53
E	52	23	12	6	12	48

Commented [MOU5]: Please discuss each finding with more updated international literature.

How is the students' learning progress level in the same item? Table 5 displays the percentage data of students' competence in explaining item 13SG-3 based on five categories of conceptual understanding (LOK, AM, MFN, MFP, and SK). In the SK category, students in group D perform better among all groups (D (53%) > E (48%) > B (31%) > C (27%) >A (19%)). Simply put, more than half students in group D are capable of explaining the item 13SG-3 compared to students in other groups. Meanwhile, in LOK, students in group A and C show higher percentage among all groups (A (36%) = C (36%) > E (23%) > B (19%) > D(13%)). In other words, more than one-third of students in group A or C is incapable of explaining the item 13SG-3 compared to students in other groups due to the limited knowledge on the item. Moroever, in AM, group B shows highest percentage among all groups (B (36%) > C (27%) > D (24%) > A (21%) > E (12%)); it signifies that more than one-fourth of students in group B are incapable of explaining item 13SG-3 compared to other groups due to the misconception on the item. Such findings indicate that the high percentage in LOK and AM category is seen as one of the reasons why the students' competence is different in explaining the same item 13SG-3. To put it another way, the students' learning progress does not develop optimally in explaining item 13SG-3 due to lack of knowledge (LOK) or misconception (AM) on the item.



Figure 2(a) Probability Category Curve of item 13SG-3 of group A, and Figure 2(b) Probability Category Curve of item 13SG-3 of group D (Category: LOK = Lack of Knowledge, AM = All-Misconception, MFN = Misconception False Negative, MFP = Misconception False Positive, SK = Scientific Knowledge)

Figure 2 illustrates the comparison of the probability category curve (PCC) of students in group A and D in item 13SG-3. The five curve shapes are the visual representation of the distribution of five categories of students' conceptual understanding. From the curves, one can identify which groups that tend to show LOK and AM category traits. It is worth noting that the curve 2(a) and 2(b) tend to be different based on the MFP curve shape, while others are relatively similar. The MFP curve of students A has a higher probability compared to that of students D; simply put, a senior high school student tends to show stronger MFP category compared to a third-year college student. The notion is supported by the finding that senior high school students are relatively incapable of providing correct reason on item 13SG-3 compared to third-year college students. On the other hand, students with low ability in group D tend to show similar curve shape of LOK, AM, and MFN with group A. This implies that both groups' conceptual understanding in the item is relatively similar. In other words, the learning progress of group D, particularly in students with low ability, has not developed

optimally despite the fact that that group D consists of third-year college students that

progressed through three years of formal education experience in university.

This echoes previous findings that the learning progress is highly dependent on the students' learning process and experience (Duschl et al., 2011; Park et al., 2017; Wilson, 2009). Learning progress is defined as a sophisticated and systematic way of thinking, in which the students will undergo gradual progress when learning a topic for a long time interval. Such a systematic way of thinking is formed by the learning practices and education experience (Emden et al., 2018). On top of that, the research findings are in line with previous studies that highlighted that students have distinctive comprehension formed by their own experience (Chi et al., 2018; Emden et al., 2018; Hoe & Subramaniam, 2016; Jin, Mikeska, Hokayem, & Mavronikolas, 2019; Rogat et al., 2011; Testa et al., 2019). Such distinctive knowledge has not been explored by evaluation or intervention through learning roadmaps that are in accordance with remedial learning (Smith et al., 2006). In spite of that, it is considered essential to conduct a further analysis that focuses on the modification of conceptual understanding category and analysis variation that is able to define the characteristics of students' alternative conception. The development procedures, as explained in the methodology, has resulted in 4TMC instrument; however, instrument development is seen as an essential continuous process (Wilson, 2009, 2012).

Based on the research findings, the study identifies several important notes on the development of the 4TMC instrument. Firstly, further analysis of the characteristic of students' response behavior is necessary to conduct regarding the item clarity and the measured concept. The findings have implied that the percentage of LOK and AM understanding category is relatively dominant and tends to increase along with the level of conceptual understanding. Hence, the development of the concept level requires taking into consideration any potential term use that might confuse the students. A further study on the identification of commonly-understood terms or concepts is therefore essential. Secondly, a

Commented [MOU6]: This is also finding please move it to the findings section.

separate analysis is required to diagnose the factors contributing to the students' lack of knowledge and misconception. Regarding that, further analysis can be conducted by applying the analysis methods developed by previous studies (Caleon & Subramaniam, 2010; Hoe & Subramaniam, 2016; Oon & Subramaniam, 2013). Thirdly, it is discovered that the concepts LG, SG, SL and LS were interpreted differently by the students. Despite being in the same conceptual understanding level, the items' difficulty level are completely different. Therefore, an evaluation on answer choices requires one to focus on the representation of understanding at the same level.

One of the features of the Rasch model is that the model facilitates one to identify any correlation between the construct map and the students' competence in ways that the students' competence can be analyzed by referring to the difference in item difficulty level. The 4TMC instrument indicates that there are students with very high ability as well as students with low ability in each group. Such a gap serves as the basis for qualitative interpretation to elaborate on the difference in students' competence. The insight is applicable in the learning process of chemistry subject. The instrument is expected to be beneficial for teachers in developing a formative test to identify the students' progress of conceptual understanding. On top of that, teachers are able to implement the instrument as a diagnostic instrument to evaluate students' conceptual understanding in providing feedback on their learning progress. Further, the teachers will be able to develop instructional strategies that are specifically designed to tackle the students' difficulty in developing an epistemological explanation regarding the concept of change of state of matter.

Conclusions

The result revealed that the integration of the 4TMC test and Rasch modeling is effective and valid to be treated as the diagnostic instrument to measure students' learning progress. Moreover, it is discovered that students in group A, B, C, D, and E, particularly those with

low ability, are hampered in developing an epistemological explanation of the concept. This blames the students' lack of certainty in their answer and reason; thus, assumed as having lack of knowledge or misconception. The low-ability students' curve shape of LOK and AM is consistent in the competence interval of less than 0.1 logit. On the other hand, the students' ability gets lower as the conceptual understanding level increases. Such finding indicates that the learning process and education experience provide a limited contribution for the students in developing a systematic way of thinking regarding the concept of change of state of matter.

Recommendations

The Based on the results of the study, there are several recommendations for researchers and teachers. For researchers, the findings of this research can be followed up to examine more in how students build their understanding gradually in explaining the concept of particles in substance form changes. The study can be conducted by developing tests that aim to evaluate and diagnose the process of student knowledge formation and development while being able to identify at the level of education where the confusion of understanding occurs. The evaluation becomes more objective, not only reviewed from the student's point of ability but can be reviewed from the teacher's ability. The model of *Rasch's* multi-faced item response pattern approach becomes one of the important parts recommended for such objectives. In this way, students' ability to develop epistemological knowledge, and their ability to significantly actualize the knowledge gained can be measured well.

On the other hand, for teachers, the results of this study along with the stages of analysis approach used can be a reference in evaluating the progress of learners' learning, as well as determining alternative thinking frameworks of students in explaining the concept of substance change. The information serves as strategic feedback in formulating instructional strategies and preparing remedial learning, especially for students who have difficulty in developing epistemological explanations of substance changes.

Limitations

The limitations of the research are primarily related to the misrepresentation of student reasoning, which may arise in its efforts to connect phenomena and concepts measured in each item. In this context, the student may not excel to explain, because of his incapableness in using his heuristic reasoning. This instrument is not equipped with items that evaluate the heuristic abilities of the student in question. However, researchers decided to record this incompetence as a misconception or vague knowledge. For further research, it is recommended that the instrument be equipped with items that measure students' emotional and heuristic reasoning according to the conceptual framework to be evaluated.

Acknowledgments

The researchers would like to express their gratitude towards the Directorate of Research and Community Service, Ministry of Research and Technology of Republic of Indonesia, for the financial support through the University Basic Research Excellence Grant Program in the Research and Community Service Office of Universitas Negeri Gorontalo, 2020.

References

- Aktan, D. C. (2013). Investigation of students' intermediate conceptual understanding levels: The case of direct current electricity concepts. *European Journal of Physics*, 34(1), 33– 43. https://doi.org/10.1088/0143-0807/34/1/33
- Arslan, H. O., Cigdemoglu, C., & Moseley, C. (2012). A three-tier diagnostic test to assess pre-service teachers' misconceptions about global warming, greenhouse effect, ozone layer depletion, and acid rain. *International Journal of Science Education*, 34(11), 1667– 1686. https://doi.org/10.1080/09500693.2012.680618
- Bond, T. G., & Fox, C. M. (2007). *Applying The Rasch Model: Fundamental Measurent in the Human Sciences* (2nd Ed.). Routledge Taylor & Francis Group

- Caleon, I. S., & Subramaniam, R. (2010). Do students know what they know and what they don't know? Using a four-tier diagnostic test to assess the nature of students' alternative conceptions. *Research in Science Education*, 40(3), 313–337. https://doi.org/10.1007/s11165-009-9122-4
- Chandrasegaran, A. L., Treagust, D. F., & Mocerino, M. (2007). The development of two tier multiple-choice diagnostic instrument for evaluating secondary school students' ability to describe and explain chemical reactions using multiple levels of representation. *Chemistry Education Research and Practice*, 8(3), 293–307
- Chi, S., Wang, Z., Luo, M., Yang, Y., & Huang, M. (2018). Student progression on chemical symbol representation abilities at different grade levels (Grades 10–12) across gender. *Chemistry Education Research and Practice*, 19(4), 1055–1064. https://doi.org/10.1039/c8rp00010g
- Claesgens, J., Scalise, K., Wilson, M., & Stacy, A. (2009). Mapping student understanding in chemistry: The perspectives of chemists. *Science Education*, 93(1), 56–85. https://doi.org/10.1002/sce.20292.
- Duncan, R. G., & Hmelo-Silver, C. E. (2009). Learning progressions: Aligning curriculum, instruction, and assessment. *Journal of Research in Science Teaching*, 46(6), 606–609. https://doi.org/10.1002/tea.20316
- Duschl, R., Maeng, S., & Sezen, A. (2011). Learning progressions and teaching sequences: A review and analysis. *Studies in Science Education*, 47(2), 123–182. https://doi.org/10.1080/03057267.2011.604476
- Emden, M., Weber, K., & Sumfleth, E. (2018). Evaluating a learning progression on "Transformation of Matter" on the lower secondary level. *Chemistry Education Research and Practice*, 19(4), 1096–1116. https://doi.org/10.1039/c8rp00137e
- Habiddin, & Page, E. M. (2019). Development and validation of a four-tier diagnostic instrument for chemical kinetics (FTDICK). *Indonesian Journal of Chemistry*, 19(3), 720–736. https://doi.org/10.22146/ijc.39218
- Hadenfeldt, J. C., Bernholt, S., Liu, X., Neumann, K., & Parchmann, I. (2013). Using ordered multiple-choice items to assess students' understanding of the structure and composition of matter. *Journal of Chemical Education*, 90(12), 1602–1608. https://doi.org/10.1021/ed3006192
- Hasan, S., Bagayoko, D., & Kelley, E. L. (1999). Misconceptions and the certainty of response index (CRI). *Physics Education*, 34(5), 294–299. https://doi.org/10.1088/0031-9120/34/5/304
- Herrmann-Abell, C. F., & Deboer, G. E. (2016). Using rasch modeling and option probability curves to diagnose students' misconceptions. *American Eduacational Research Assossiation*, 1–12
- Hoe, K. Y., & Subramaniam, R. (2016). On the prevalence of alternative conceptions on acidbase chemistry among secondary students: Insights from cognitive and confidence measures. *Chemistry Education Research and Practice*, 17(2), 263–282.

https://doi.org/10.1039/c5rp00146c

- Jin, H., Mikeska, J. N., Hokayem, H., & Mavronikolas, E. (2019). Toward coherence in curriculum, instruction, and assessment: A review of learning progression literature. *Science Education*, 103(5), 1206–1234. https://doi.org/10.1002/sce.21525
- Klassen, S. (2006). Contextual assessment in science education: Background, issues, and policy. *Science Education*, 90(5), 820–851. https://doi.org/10.1002/sce.20150
- Laliyo, Botutihe, & Panigoro. (2019). The development of two-tier instrument based on distractor to assess conceptual understanding level and student misconceptions in explaining redox reactions. *International Journal of Learning, Teaching and Educational Research*, 18(9), 216–237. https://doi.org/10.26803/ijlter.18.9.12
- Linacre, J. M. (2012). A user's guide to WINSTEPS ® MINISTEP Rasch-model computer program: Program manual 3.75.0. https://doi.org/ISBN 0-941938-03-4
- Linacre, J. M. (2020). A User's Guide to WINSTEPS ® MINISTEP Rasch-Model Computer Programs Program Manual 4.5.1. https://doi.org/ISBN 0-941938-03-4
- Ling Lee, W., Chinna, K., & Sumintono, B. (2020). Psychometrics assessment of HeartQoL questionnaire: A Rasch analysis. *European Journal of Preventive Cardiology*, 1–6. https://doi.org/10.1177/2047487320902322
- Liu, X. (2012). Developing measurement instruments for science education research. In B. Fraser, K. G. Tobin, & C. J. McRobbie (Eds.), Second international handbook of science education (pp. 651–665). Springer Netherlands
- Lu, S., & Bi, H. (2016). Development of a measurement instrument to assess students' electrolyte conceptual understanding. *Chemistry Education Research and Practice*, 17(4), 1030–1040. https://doi.org/10.1039/c6rp00137h
- Morell, L., Collier, T., Black, P., & Wilson, M. (2017). A construct-modeling approach to develop a learning progression of how students understand the structure of matter. *Journal of Research in Science Teaching*, 54(8), 1024–1048. https://doi.org/10.1002/tea.21397
- Neumann, K., Viering, T., Boone, W. J., & Fischer, H. E. (2013). Towards a learning progression of energy. *Journal of Research in Science Teaching*, 50(2), 162–188. https://doi.org/10.1002/tea.21061
- Park, M., Liu, X., & Waight, N. (2017). Development of the connected chemistry as formative assessment pedagogy for high school chemistry teaching. *Journal of Chemical Education*, 94(3), 273–281. https://doi.org/10.1021/acs.jchemed.6b00299
- Peterson, R. F., Treagust, D. F., & Garnett, P. (1989). Development and application of a diagnostic instrument to evaluate grade-11 and -12 students' concepts of covalent bonding and structure following a course of instruction. *Journal of Research in Science Teaching*, 26(4), 301–314. https://doi.org/10.1002/tea.3660260404
- Rogat, A., Anderson, C., Foster, J., Goldberg, F., Hicks, J., Kanter, D., ... Wiser, M. (2011). Developing learning progression in support of the new science standards: A RAPID

workshop series. (4), 163. https://doi.org/10.12698/cpre.2011.lprapid

- Sadler, P. M. (1999). The relevance of multiple-choice testing in assessing science understanding. In J. J. Mintzes, J. H. Wandersee, & J. D. Novak (Eds.), Assessing science understanding: A human constructivist view (pp. 251–274). Zaccheus Onumba Dibiaezue Memorial Libraries. https://zodml.org/sites/default/files/%5BJoel_J._Mintzes%2C_James_H._Wandersee%2 C_Joseph_D._No_0.pdf
- Smith, C. L., Wiser, M., Anderson, C. W., & Krajcik, J. (2006). Implications of research on children's learning for standards and assessment: A proposed learning progression for matter and the atomic-molecular theory. *Measurement: Interdisciplinary Research & Perspective*, 4(1–2), 1–98. https://doi.org/10.1080/15366367.2006.9678570
- Sumintono, B., & Widhiarso, W. (2014). Aplikasi model Rasch untuk penelitian ilmu-ilmu sosial [Application of Rasch model in social science research]. Trim Komunikata. https://www.researchgate.net/publication/268688933%0AAplikasi
- Testa, I., Capasso, G., Colantonio, A., Galano, S., Marzoli, I., Scotti di Uccio, U., ... Zappia, A. (2019). Development and validation of a university students' progression in learning quantum mechanics through exploratory factor analysis and Rasch analysis. *International Journal of Science Education*, 41(3), 388–417. https://doi.org/10.1080/09500693.2018.1556414
- Treagust, D. F. (1988). Development and use of diagnostic tests to evaluate students' misconceptions in science. *International Journal of Science Education*, 10(2), 159–169. https://doi.org/10.1080/0950069880100204
- Tyson, L., Treagust, D. F., & Bucat, R. B. (1999). The complexity of teaching and learning chemical equilibrium. *Journal of Chemical Education*, 76(2–4), 554–558. https://doi.org/10.1021/ed077p1560.1
- Wilson, M. (2005). Constructing measures: an item response modeling approach. Lawrence Erlbaum Associates, Inc. https://doi.org/10.4324/9781410611697
- Wilson, M. (2008). Cognitive diagnosis using item response models. Zeitschrift Für Psychologie / Journal of Psychology, 216(2), 74–88. https://doi.org/10.1027/0044-3409.216.2.74
- Wilson, M. (2009). Measuring progressions: Assessment structures underlying a learning progression. Journal of Research in Science Teaching, 46(6), 716–730. https://doi.org/10.1002/tea.20318
- Wilson, M. (2012). Responding to a challenge that learning progressions pose to measurement practice. In A. C. Alonzo & A. W. Gotwals (Eds.), *Learning progression* in science (pp. 317–344). https://doi.org/10.1007/978-94-6091-824-7



European Journal of Educational Research

ISSN: 2165-8714

http://www.eu-jer.com/

Review Form					
Manuscript ID:	Manuscript ID: EU-JER_ID# Date:				
Manuscript Title:Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress					
(Ma	ABOUT MANUSCRIPT rk with "X" one of the options)	Accept	Weak	Refuse	Not Available
Language is clear a	nd correct	x			
Literature is well w	ritten	X			
References are cite	d as directed by APA		x		
The research topic	is significant to the field	X			
The article is compl	ete, well organized and clearly written		x		
Research design an	d method is appropriate	X			
Analyses are appro	priate to the research question	x			
Results are clearly p	presented		x		
A reasonable discus	ssion of the results is presented		x		
Conclusions are cle	arly stated	X			
Recommendations are clearly stated x					
	GENERAL REMARKS AND RECOMME	NDATIONS TO TH	IE AUTHOR		
 The manuscript is related to the implementation of four-tier instruments based on the Rasch Model in evaluating students' learning progress about states of matter. It has some structural deficits. The following recommendations are presented: 1- Please double-check that all citations are in the text and the references are fitting to APA 7. 2- Why did you ignore gas-solid(GS) and gas-liquid(GL)? 3- There are findings in the discussion section please move it to the findings section. 4- Please discuss each finding with more updated international literature. 					
THE DECISION (Mark with "X" one of the options)					
Accepted: Correction	on not required				
Accepted: Minor co	prrection required				
Conditionally Accepted: Major Correction Required (Need second review after corrections) x					x
Refused					
Reviewer Code: R2611 (The name of referee is hidden because of blind review)					



EuropeanJournal of EducationalResearch

ISSN: 2165-8714

http://www.eu-jer.com/

Review Form						
Manuscript ID:	EU-JER_ID#2011240749 Date: 02/07/2021					
ManuscriptTitle:	le: Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress					
ABOUT MANUSCRIPT (Mark with "X" one of the options) Accept Weak Refu					Not Available	
Language is clear ar		X				
Literature is well wi		X				
References are cite	X					
The research topic		X				
The article is compl	ete, well organized and clearly writ	en	X			
Research design an	d method is appropriate		X			
Analyses are appro	priate to the research question		X			
Results are clearly p	presented		X			
A reasonable discus		X				
Conclusions are cle	X					
Recommendations	are clearly stated		X			
	GENERAL REMARKS AND REC	OMMENDATIONS TO TH	IE AUTHOR			

1 Abstract

1.1Statistics should not be reported in the abstract.

- 2 Introduction
- 2.1 I will suggest dividing the introduction into several parts (e.g., students' understanding of matters, learning progressions, 4TMC) to make it easier for readers to understand.
- 2.2 The literature review cannot well elicit the research motivation and problems of this study. As a reader, I am interested in why an instrument in the-4TMC format was used to diagnose and investigate the progressions of students' conceptual understanding of the change of matter states.
- 2.3 The notion of "learning progress" should be replaced by "learning progressions."
- 2.4 The authors stated that "Four-tier instruments are used in studies discussing topics such as ..." and "that studies on chemistry topic which employ four-tiers instruments tend to focus on describing alternative conception." However, as I know, the greatest strength of 4TMC is that it can provide cognitive and confidence measures (e.g., Sreenivasulu& Subramaniam, 2014) at the same time. The cognitive measures can be interpreted as students' conceptual understanding or students' proficiency in practices, which are the focus of many previous LP studies. I appreciate the author's contribution to introducing 4TMC in LP studies; however, it is not clear how does the 4TMC test differ from other instruments used in previous LP studies.
- 3 Methodology
- 3.1 The study involved four concepts of change of state of matter: liquid-gas (LG), solid-liquid (SL), solid-gas (SG), and liquidsolid (LS). Why do the other two types of change of state of matter (i.e., gas-liquid (GL) and gas-solid (GS)) were not touched in the current study?
- 3.2 What is the basis of the five levels of conceptual understanding? Do the hypothesized based on previous studies (e.g., Hadenfeldt et al., 2013 or their research published in JRST Hadenfeldt, J. C., Neumann, K., Bernholt, S., Liu, X., &Parchmann, I. (2016). Students' progression in understanding the matter concept. Journal of Research in Science Teaching. http://onlinelibrary.wiley.com/doi/10.1002/tea.21312/pdf)? The author should learn the process of how Hadenfeldt et al. (2013, 2016) develop the levels of conceptual understanding.
- 3.3 It would be nice to provide examples for item design and explain the process of items and their options to match the assumed levels of understanding.
- 3.4 Why do you choose the participants comprised of both high school students and college students? Does the level of conceptual understanding match the course they are studying? More background information about the participants should be provided.



EuropeanJournal of EducationalResearch

ISSN: 2165-8714

http://www.eu-jer.com/

- 3.4 The author stated that "The involvement of Rasch model measurement lay on the assumption that the item difficulty level is dependent on the students' answer, and that the students' understanding is dependent on the estimation of item difficulty (Linacre, 2012)." This statement may be incorrect. In Rasch modeling, person and item measures are independent.
- 3.5 The PCM model should be written in a separate line. The statistical package used for Rasch analysis and ANOVA were not identified in the manuscript.

4. Results and Discussion

- 4.1 The evidence for the unidimensionality was inadequate. The eigenvalue of the 1st contrast produced by the PCA on Rasch residuals should be provided.
- 4.2 The sequence of results of the Rasch analysis should be modified. See an example in Yang, He, and Liu (2017). Yang, Y., He, P., & Liu, X. (2017). Validation of an Instrument for Measuring Students' Understanding of Interdisciplinary Science in Grades 4-8 over Multiple Semesters: A Rasch Measurement Study. International Journal of Science and Mathematics Education, 1–16.
- 4.3 The measures and fit statistics of items and the Wright map should be provided. The match between item measures and person measures can establish evidence for your instrument.
- 4.4 The authors stated that "Moreover, the findings also suggest that the item difficulty level (LG, SL, SG, and LS; particularly SL (melting) and LS (freezing)) do not match the level in the construct map." I will suggest the authors modify the items (or the hypothesized levels of conceptual understanding) before subsequent inferential statistics.
- 4.5 The representation of the data in figure 1 is quite different from that in the current literature. The standard deviation was always plotted as an error line on the mean value.
- 4.6 The item characteristic curves plotted in Figure 2 suggested that the current study's scoring method may be questionable.

Overall comment:

For me, it becomes not apparent why we need a new study on Rasch analysis of a four-tier instrument for assessing students' progression on the change of state matter. I will suggest the author focus more on the instrument's content (i.e., the change of state matter) or the methodology of scoring issues faced by 4TMC. In my opinion, the latter issue may be more significant for the science education community. Students' response to a 4TMC item can be classified into five categories, i.e., Lack of Knowledge, All-Misconception, Misconception False Negative, Misconception False Positive, Scientific Knowledge. However, it is unclear to what extent do these categories can be distinguished based on students' latent ability estimated by the Rasch model. Further studies should be conducted to reevaluate the scoring methods of 4TMC.

THE DECISION (Mark with "X" one of the options)				
Accepted: Correction not required				
Accepted: Minor correction required				
Conditionally Accepted: Major Correction Required (Need second review after corrections)	X			
Refused				
Reviewer Code: R2615 (The name of referee is hidden because of blind review)				



European Journal of Educational Research

ISSN: 2165-8714

http://www.eu-jer.com/

Review Form						
Manuscript ID:	EU-JER_ID# 2011240749	Dat	t <mark>e:</mark> Fel	oruary 2, 20	20	
Manuscript Title:						
(Ma	ABOUT MANUSCRIPT (Mark with "X" one of the options)			Weak	Refuse	Not Available
Language is clear a	nd correct			X		
Literature is well w	ritten		X			
References are cite	d as directed by APA			X		
The research topic	is significant to the field		X			
The article is comp	lete, well organized and clearly writt	en	X			
Research design an	d method is appropriate			X		
Analyses are appro	priate to the research question			X		
Results are clearly	presented			X		
A reasonable discus	ssion of the results is presented			X		
Conclusions are cle	arly stated			X		
Recommendations	are clearly stated			X		
	GENERAL REMARKS AND RECO	OMMENDATIO	ONS TO TH	IE AUTHOR		
 In the content section of the article, correct it according to the suggestions and input listed in the article script 						
THE DECISION (Mark with "X" one of the options)						
Accepted: Correction not required					V	
Conditionally Acco	nted: Major Correction Pequired (N	lead second r	oviow ofto	r correction		^
				i correction		
Re	eviewer Code: R2614 (The name of	referee is hide	den becau	se of blind r	eview)	



European Journal of Educational Research

ISSN: 2165-8714

http://www.eu-jer.com/

Review Form					
Manuscript ID:	EU-JER_ID# 2011240749 Date: February 7, 2021				
Manuscript Title:	Implementation of Four-tier Instruments Learning Progress	Based on the Ra	asch Model i	n Evaluatin	g Students'
(Mai	ABOUT MANUSCRIPT 'k with "X" one of the options)	Accept	Weak	Refuse	Not Available
Language is clear an	nd correct	X			
Literature is well wr	itten	X			
References are cited	d as directed by APA		X		
The research topic i	s significant to the field	X			
The article is compl	ete, well organized and clearly written		X		
Research design and	d method is appropriate	X			
Analyses are approp	priate to the research question	X			
Results are clearly p	resented	X			
A reasonable discus	sion of the results is presented		X		
Conclusions are clea	arly stated		X		
Recommendations	are clearly stated	X			
	GENERAL REMARKS AND RECOMMEN	IDATIONS TO TH	IE AUTHOR		
Much needed impro Much needed Appe	ovement in Abstract, introduction, methods ndix at the end of the paper by attaching th	and discussions e 4TMC instrume	ent that has	been develo	ped.
Discuss each step of	^f your study.				
Discuss why the dia	gnostic instrument test can measure misco	nceptions.			
Include the results of previous study on the development of misconception instruments, then discuss the privilege of this instrument (4TMC) compared with the previous instrument. Also, discuss the benefits of this study for education.					e privilege 'y for
The author should have revised all items that gave from the reviewed results. You can see the review result based on the manuscripts revision.					
THE DECISION (Mark with "X" one of the options)					
Accepted: Correction not required					
Accepted: Minor correction required X					X
Conditionally Accepted: Major Correction Required (Need second review after corrections)					
Refused					
Re	viewer Code: R2612 (The name of referee	is hidden becau	se of blind r	eview)	

	CORRECTION REPORT				
No	Reviewer Code	Reviews	Corrections made by the author		
1					
2					
3					
4					
5					
6					
7					

8		
9		
10		
11		
12		
13		
14		

15		
16		
17		
18		
19		
20		
21		

22		
23		
24		
25		
26		
~~		
27		
28		

29		
30		
31		
32		
33		
34		
35		

36		
37		
38		
39		
40		
41		
42		

43		
44		
45		
46		
47		
48		
49		

50		
51		
52		
53		
54		
55		

Implementation of Four-tier Multiple-choice Instruments Based on the Partial Credit Model in Evaluating Students' Learning Progress

Abstract: One of the issues that hinder the students' learning progress is the inability to construct an epistemological explanation of a scientific phenomenon. Four-tier multiple-choice (hereinafter, 4TMC) instrument and Partial-Credit Model model were employed to elaborate on the diagnosis process of the aforementioned problem. The purpose of this study was to developing and implementation four-tier multiple-choice instrument with Partial-Credit Model to evaluate students' learning progress in explaining the concept of change of state of matter. This research used a development research referring to the test development model from Wilson. The data were obtained through development and validation techniques on 20 4TMC items tested to 427 students. On each item, the study applied diagnostic-summative assessment and certainty response index. The students' conceptual understanding level was categorized based on the combination their answer choices; the measurement generated Partial-Credit Model for 1 parameter logistic (IPL) data. Analysis of differences based on class level of students using Analysis of Variants (One-way ANOVA). This study succeeded in developing 20 valid and reliable 4TMC instruments. The result revealed that the integration of 4TMC test and Partial-Credit Model was effective to be treated as the instrument to measure students' learning progress. One-way ANOVA test indicates a difference among the students' competence based on the academic level. On top of that, it was discovered that low-ability students see very slow progress due to the lack of knowledge as well as a misconception in explaining the Concept of Change of State of Matter. All in all, the research regarded that the diagnostic information was necessary for teachers in prospective development of learning strategies and evaluation of science learning.

Keywords: Learning progress. four-tier, change of state of matter, Partial-Credit Model.

Introduction

Central to the notion of science learning is the development of students' scientific understanding of basic concepts of sciences (Hadenfeldt et al., 2013), particularly, change of state of matter (Emden et al., 2018). Aside from the issue, several studies have also highlighted the students' inability to provide an epistemological explanation of basic concepts of sciences (Chi et al., 2018). Efforts to solve the issues, however, have shown little progress,

as the students might have more complex perceptions regarding the alternative concept they understand (Morell et al., 2017).

Education practitioners have recommended the utilization of learning progress concept as the instructional method to provide guidance and direction and to adjust the curriculum with the learning process and assessment (Claesgens *et al.*, 2009; Duncan & Hmelo-Silver, 2009; Rogat *et al.*, 2011). Learning progress is defined as a sophisticated and systematic way of thinking. This method is applicable for a learning process, in which the students will undergo gradual progress when learning a topic in a long duration. Its effectiveness is highly dependent on the learning process and the students' learning experience (Duschl et al., 2011). The concept involves certain sets of gradual levels that represent conceptual understanding, ranging from low level up to comprehensive level.

The notion of learning progress is highly distinctive to each student and is dependent to one's learning experience (Rogat et al., 2011); therefore, there is no learning roadmap that is suitable for all kinds of students (Smith et al., 2006). Each student constructs one's understanding in a different way; moreover, the construction process is varied depending on the students' conceptual understanding level (Aktan, 2013). This is to say that each student undergoes a different rate of learning progress, understanding level, and knowledge construction. Simply put, the development of scientific comprehension among students is not linear (Neumann et al., 2013). Therefore, this study regards each level of students' conceptual understanding for more advanced level of understanding (Hadenfeldt et al., 2013). A student who faces difficulty in a certain level of understanding will see a lack of progress to a more advanced level. This in turn hinders the student's ability to construct an epistemological explanation on the basic concepts of science. Within this context, the learning progress is treated as the method to evaluate students' conceptual understanding. The diagnostic information generated is reliable to be treated as a reference for

the teachers in developing accurate and valid instructional components to guide the students to progress to the next level.

Among the diagnostic instruments that are considered applicable is the **four-tier multiplechoice (4TMC)** instrument. It is the development of two-tier multiple-choice test recommended by Treagust (1988) and Chandrasegaran et al., (2007). The use of two-tier instrument is familiar in identifying students' understanding in select topics such as electrochemistry (Lu & Bi, 2016), covalent bond (Peterson et al., 1989), and chemical equilibrium (Tyson et al., 1999). Despite its reputation in academia, the two-tier test has raised criticism due to its sole focus on the facts and negligence towards students' understanding (Klassen, 2006). Therefore, several experts propose the renewed version of the test by adding distractor answer choices to strengthen the diagnostic value of the items (Herrmann-Abell & DeBoer, 2011; Sadler, 1998). In addition, some have highlighted the test's weakness in cases where students' tended to pick the answer choice and the reasoning randomly. This illustrates that the students were uncertain and possessed several misconceptions in the first tier question. In such cases, teachers faced difficulty in differentiating between guessed answers and misconceptions (Habiddin & Page, 2019; Hasan et al., 1999).

The criticism laid against the model has sparked the innovation of three-tier and four-tiers instruments. Both instruments feature two multi-level questions, also similar with two-tier test. In the three-tier test, however, the measurement of students' certainty level is conducted simultaneously in both first and second-tier questions; in the meantime, the measurement is conducted separately in the first two tiers (Caleon & Subramaniam, 2010). The value of students' certainty rate ranges from one (very uncertain) to five (very certain).

Three-tier test lacks validity in measuring the students' certainty rate regarding both the answer choice and the reasoning, whether or not the value of certainty rate refers only to the answer choice, to the reasoning, or both. Such weakness will in turn obstructs the evaluation

and classification process of students' responses (Arslan et al., 2012). In the four-tier instrument, the measurement of certainty rate also involves the answer choice in the first tier and the reasoning in the third tier (Arslan et al., 2012; Loh et al., 2014). Regarding this feature, four-tier test is considered more accurate than the three-tier test. Students who pick wrong answer choices with high certainty indicate that they have a very high misconception on the measured item (Hoe & Subramaniam, 2016).

Four-tier instruments are used in studies discussing topics such as physics education (Caleon & Subramaniam, 2010), chemical thermodynamics (Sreenivasulu & Subramaniam, 2013), transition metal (Sreenivasulu & Subramaniam, 2014), acid-base reaction (Hoe & Subramaniam, 2016), and chemical kinetics (Habiddin & Page, 2019). However, it is worth noticing that studies on chemistry topic which employ four-tiers instruments tend to focus on describing alternative conception. To put it another way, the higher the certainty rate is, the stronger the students' alternative conception will be. Despite its potentials, the scholarly discussion has overlooked the implementation of a four-tier diagnostic instrument to measure students' learning progress. Therefore, further analysis is essential on the application of 4TMC test in several domains analyzes by Partial-Credit Model approach.

The use of Partial-Credit Model has been introduced since the 2000s in the science education research; it features the instrument that integrates diagnostic assessment and summative assessment (Liu, 2012; Wei et al., 2012). On top of that, the diagnostic assessment approach is introduced to conduct an in-depth analysis of the construction process of students' conceptual understanding (Claesgens et al., 2009; Hadenfeldt et al., 2013; Lu & Bi, 2016). This study employs 4MTC and Partial-Credit Model as a diagnostic tool to evaluate students' learning progress in explaining the change of state of matter, besides focusing on the Concept of Change of State of Matter, this research employs in-depth analysis using Item Response Theory, namely Partial Credit Model.

One of the features of the Partial-Credit Model is that the model facilitates one to identify any correlation between the construct map and the students' competence in ways that the students' competence can be analyzed by referring to the difference in item difficulty level. The 4TMC instrument indicates that there are students with very high ability as well as students with low ability in each group. Such a gap serves as the basis for qualitative interpretation to elaborate on the difference in students' competence. The insight is applicable in the learning process of chemistry subject. The instrument is expected to be beneficial for teachers in developing a formative test to identify the students' progress of conceptual understanding. On top of that, teachers are able to implement the instrument as a diagnostic instrument to evaluate students' conceptual understanding in providing feedback on their learning progress. Further, the teachers will be able to develop instructional strategies that are specifically designed to tackle the students' difficulty in developing an epistemological explanation regarding the concept of change of state of matter. The study focus revolves around three research questions: 1) What is the quality of the developed 4TMC instrument based on the Partial-Credit Model?. 2). How is the effectiveness of 4TMC instrument to evaluate the students' learning progress in explaining concepts of change of state of matter. 3) How is the learning progress in students ranging from the senior high school level up to the senior (fourth) year of college in explaining the concepts?

Methodology

Development Model

This research used a development research referring to the test development model from Wilson. Wilson (2005, 2008) introduces four steps of measurement instrument development in figure 1.



Figure 1. Measurement instrument development

This recommendation is proven valid to be implemented in developing measurement instrument for different construct variables (Barbera, 2013; Chi et al., 2018; Hadenfeldt et al., 2013; Laliyo, Botutihe, & Panigoro, 2019; Lu & Bi, 2016; Wei et al., 2012; Wilson, 2009; Wind, Tsai, Grajeda, & Bergin, 2018). The present study also included two questions related to certainty rate (Arslan et al., 2012; Habiddin & Page, 2019: Hasan et al., 1999). The obtained data were analyzed by Partial Credit Model (PCM) approach by WINSTEPS version

4.5.3 software.

Construct Map: Determining Level of Understanding

The first step was to develop the construct of measured variables. The study involved four concepts of change of state of matter: liquid-gas (LG), solid-liquid (SL), solid-gas (SG), and liquid-solid (LS). Gas-liquid (GL) and Gas-Solid (GS) materials were not included in this study as they are included in the basic level of knowledge. The change of a substance from gas to solid (GS) is known as freezing, while from gas to liquid (GL) is called condensing. These two types of changes in the form of substances are very easy to answer by students at a higher level since the materials have always been presented in textbooks, from high school to university students, on the topic of changes in the form of the substance. These concepts were implemented in a gradual manner through five levels of conceptual understanding (Table 1). Such method functions as the pathway of conceptual development that involves learning
objectives from the lowest to the highest level of conceptual understanding (Duncan & Hmelo-Silver, 2009; Löfgren & Helldén, 2009; Hadenfeldt et al., 2013; Rogat et al., 2011). In other words, the set of levels, as mentioned previously, was adjusted to the students' needs so as to develop their conceptual understanding. This took into account that each student might progress on different and non-linear development of conceptual understanding; therefore, the levels, as illustrated in Table 1, was considered valid to illustrate the ideal conceptual development pathway (Neumann et al., 2013).

Table 1. Level of Conceptual Understanding in Explaining Concept of Change of State of Matter

	Concentual Understanding Level	Change of State of Matter/Item				
	Conceptual Onderstanding Level	LG	SL	SG	LS	
5	Submicroscopic diagram of change of	5LG-5	10SL-5	15SG-5	20LS-5	
	state of matter					
4	Correlation between state of matter and	4LG-4	9SL-4	14SG-4	19LS-4	
	the process of change of state of matter					
3	Process of change of state of matter	3LG-3	8SL-3	13SG-3	18LS-3	
2	Concept of state of matter	2LG-2	7SL-2	12SG-2	17LS-2	
1	Factual phenomenon of state of matter	1LG-1	6SL-1	11SG-1	16LS-1	
•		10 111	та	1 1	`	

Description: (LG = liquid-gas, SL = solid-liquid, SG = solid-gas, LS = liquid-gas)

Item Design and Assessment Scheme

The second phase involved an item design. In the 4TMC instrument, all the items consisted of four-tier multiple-choices. To put it another way, each item contains four questions that combine between diagnostic-summative test (Hoe & Subramaniam, 2016; Lu & Bi, 2016; Treagust, 1988) with certainty response index (hereinafter, CRI) test (Arslan et al., 2012; Hasan et al., 1999). The first-tier questions (Q1) aimed to identify whether or not the students understand the content. Moreover, questions in the second tier (Q2) were employed to clarify the students' certainty regarding their answers in the Q1. Third-tier questions (Q3) functioned to diagnose the students' reasoning regarding their answers in the Q1. Further, questions in the second tier (Q4) were employed to clarify the students' certainty regarding their answers in the Q3. Q1 and Q3 questions in each item involved five answer choices; one among them was the correct answer, while three were the distractor, and another answer choice was open-

ended answer choice. This open-ended option allows the students to decide the answer by themselves, should they find no correct answer as in accordance with their conceptual understanding. In the meantime, the Q2 and Q4 questions involved two close-ended answer choices; the first choice was for those who are uncertain of their answer, and the second choice was for the students who are very certain of their answer (Arslan et al., 2012). The distractor choices were employed in Q1 and Q3 questions to validate the diagnostic strength of the questions (Herrmann-Abell & DeBoer, 2011; Sadler, 1998). Therefore, in the Q1 and Q3 tiers, the students would have only 0.20 or 20 percent probability of choosing the correct answer. The item Category of Grade of Students' Conceptual Understanding in Table 2.

	Que	estions			Poting		
Q1	Q2	Q3	Q4		Conceptual Onderstanding Categor	. y	Kaung
Correct	Certain	Correct	Certain	CCCC	Scientific Knowledge	SK	<mark>5</mark>
Correct	Certain	Incorrect	Certain	CCIC	Misconception False Positive	MFP	<mark>4</mark>
Incorrect	Certain	Correct	Certain	ICCC	Misconception False Negative	MFN	<mark>3</mark>
Incorrect	Certain	Incorrect	Certain	ICIC	All-Misconception	AM	2
Correct	Certain	Correct	Uncertain	CCCU	Lack of Knowledge	LOK	1
Correct	Certain	Incorrect	Uncertain	CCIU			
Correct	Uncertain	Correct	Certain	CUCC			
Correct	Uncertain	Correct	Uncertain	CUCU			
Correct	Uncertain	Incorrect	Certain	CUIC			
Correct	Uncertain	Incorrect	Uncertain	CUIU			
Incorrect	Certain	Correct	Uncertain	ICCU			
Incorrect	Certain	Incorrect	Uncertain	ICIU			
Incorrect	Uncertain	Correct	Certain	IUCC			
Incorrect	Uncertain	Correct	Uncertain	IUCU			
Incorrect	Uncertain	Incorrect	Uncertain	IUIU			

Table 2 Category of Grade of Students' Conceptual Understanding^{*)}

("Hasan, Bagayoko and Kelley, 1999; Arslan, Cigdemoglu and Moseley, 2012; Habiddin and Page, 2019) As an illustration, in the item 13SG-3, a student picks A in Q1, "very certain" in Q2, A in Q3, and "very certain" in Q4; the combination of the student's answers is ICIC. The result illustrates that the student's answer is incorrect in the Q1 and is very certain of one's error (Q2). Moreover, s/he also provides an incorrect answer in Q3 and is very certain of one's incorrect answer in Q3 (Q4). This indicates that in the item 13SG-3, the student is categorized to have all-misconception understanding (AM). In the Conceptual Understanding Category table, the category is included in fourth grade. Incorporation of the students' answer combinations in each item into the category and grade of students' understanding would result in specific data that are in accordance with the Partial-Credit Model.

Outcome Space and Data Collection

The third step involved the design of the outcome space of the correlation between items and construct maps (Bond & Fox, 2007; Wilson, 2009). The item validation was conducted independently by three expert validators to evaluate the extent of correlation between answer choices in Q1-Q3 in each item and the level of students' conceptual understanding. The validators were asked to clarify that the questions are easy to understand and the students' lack of linguistic competence would not hinder them from providing the right answer. The validators also required to ensure that the questions are in accordance with the syllabus, particularly with the students' conceptual understanding as based on the construct map. The questions in each item were also validated in several aspects, such as: ambiguity, time allocation, directiveness towards a particular answer, and subjective or emotional expression. Fleiss κ measure was employed to acquire information on the validators' approval. From the measure, it was generated that the κ value = 0.97, indicating that the three validators agreed that the 4TMC items were valid in correlating between the answer choices and the students' conceptual understanding.

The next step was to acquire data based on the measurement instrument. The instrument was tested to 427 students in Gorontalo, Indonesia using cluster random sampling technique. The students comprised 171 (40.05%) senior high school students (or students A), 83 (19.44%) university freshmen majoring chemistry education (or students B), 66 (15.45%) second-year university students majoring chemistry education (or students C), 55 (12.88%) third-year university student majoring chemistry education (or students D), and 52 (12.18%) fourth-year university students majoring chemistry education (or students E). Based on gender, the female participants comprised 369 participants (86.41%), and the male counterparts consisted of 58 participants (13,58%). The participants were given no particular educational treatments and had stated their voluntary consent to participate in the research.

Partial-Credit Model Measurement and Data Analysis

The fourth step was to conduct the Partial-Credit Model measurement. This step was implemented to define the correlation between the score generated and the students' conceptual understanding level as elaborated within the construct map. The involvement of Partial-Credit Model measurement lay on the assumption that the item difficulty level is dependent on the students' answer, and that the students' understanding is dependent on the estimation of item difficulty (Linacre, 2012).

Partial credit model (PCM) was employed to evaluate the learning progress through structured questions; this took into account that the instrument items involved gradual and structured questions (Bond and Fox, 2007; Masters, 1982; Sumintono and Widhiarso, 2015; Wilson, 2009). The model was stated into the following formula: $ln[P_nik/(1 - P_nik)] - B_n - D_ik$, in which P_{nik} refers to the probability of student n with B_n ability to pick correct response in the level k of item i; while D_{ik} refers to the difficulty level k of item i, or the threshold point for the test taker who scores k, not k -1. Analysis of differences based on class level of students using One-way ANOVA.

Results and Discussion[A2]

Results

The developed 4TMC instrument adapts the two-level instrument model by Treagust (1988), combined with the CRI theory by Arslan (2012). The function of CRI (certainty response index) is to ensure that students' choice of answers in Q1 and Q3 are the answers that they believe in. This is called diagnostic because it investigates the level of student error in stages, including the ability of students to understand and to use their understanding in explaining the reasons for their choice of answers. Thus, measurement is conducted both at the level of knowledge and reasoning.

The item design referred to the basic criteria to ensure that the students would be able to identify logical reason in Q3 as based on their answer in Q1; moreover, the item design also aimed to clarify the students' certainty of their answers through Q2 and Q4 questions. The 4TMC instrument also allows the students to state their certainty level of Q1 and Q3 answer choices separately. Students with correct understanding regarding the concept of change of state of matter (Q1) and its reasoning (Q3) will pick the "very certain" answer in the Q2 and Q4. If the students are uncertain of their answer regarding the content (Q1) but are certain of the reasoning (Q3), this suggests that the students are able to comprehend the concept/theory but unable to implement such concepts. This study views that it is beneficial to explore potential combinations of Q1/Q3 answer choices and Q2/Q4 certainty rate implementation to provide in-depth elaboration on students' understanding of certain concepts (Habiddin, 2019).

The item design is illustrated in Figure 2.



Each combination of students' answers in each item was categorized based on the assessment scheme. Every correct response in Q1 and Q3 is labeled with C (correct) code, and wrong answers were labeled with I (incorrect) code. Moreover, in the Q2 and Q4 tier, "very certain" answers are labeled with C (certain) code, and "uncertain" answers were labeled with U (uncertain) code. Therefore, combination of answer choice in each item was generated and written in sequence based on the questions in the Q1, Q2, Q3, and Q4. Such combination was treated as a reference in determining category and grade of students' conceptual understanding in each item, as shown in Table 2. The students' conceptual understanding category was adapted from findings reported by Hasan *et al.*, (1999)[A3], Arslan *et al.*, (2012), and Habiddin (2019).

Each combination of students' answer in each item was classified based on five categories of conceptual understanding. The first category is scientific knowledge (SK) with a grade of five; it illustrates that the students possess knowledge that is scientifically correct. The second category is misconception false positive (MFP) (with a grade of four), illustrating that the students have a correct claim of understanding, but they are unable to explain the claim. The third category is misconception false negative (MFN) (with a grade of three), illustrating that the students do not have correct claims of knowledge, but they are able to explain the claim. This category is considered negative because it is possible that the answers provided are guessed answer that is coincidentally correct. The fourth category is all-misconception (AM) with a grade of four; this category signifies that the students are very certain of their incorrect knowledge. Lastly, the fifth category is lack of knowledge (LOK) with a grade of five, signifying that the students lack knowledge in a particular item. Such categories were determined by the students' certainty level in Q2 or Q4. As an instance, one of the possible answer combinations is as follows: CCCU. This combination illustrates that the student is correct in Q1, Q2, and Q3, but is uncertain in Q4; the condition signifies that the student's understanding is ambiguous, hesitant, and is not based on appropriate scientific knowledge.

Effectiveness of Measurement Instruments

Unidimensionality Unidimensionality is an essential indicator to evaluate the 4TMC instrument's ability to measure students' capability of explaining the concept of change of state of matter. This indicator is measured by Principal Component Analysis of the residuals to estimate the extent of variance to which the instrument is able to measure what it is supposed to measure (Sumintono & Widhiarso, 2014). As displayed in the figure 3, the result of raw variance explained by measures of data is 38.9%, the number almost approaches the expectation value of 39.2%. The numbers indicate that the minimum unidimensionality requirements of 20% are achieved, and simultaneously, the limit of PCM unidimension is met (approaching 40%) (Linacre, 2012; Ling Lee et al., 2020). Moreover, the instrument's unexplained variance values are below 7% and considered as ideal (not exceeding 15%), signifying that the item independence rate in instrument falls into "good" category.

		Em	pirical		Modeled
Total raw variance in observations	=	32.7	100.0%		100.0%
Raw variance explained by measures	=	12.7	38.9%		39.2%
Raw variance explained by persons	=	3.8	11.7%		11.8%
Raw Variance explained by items	=	8.9	27.2%		27.5%
Raw unexplained variance (total)	=	20.0	61.1%	100.0%	60.8%
Unexplned variance in 1st contrast	=	2.0	6.2%	10.2%	
Unexplned variance in 2nd contrast	=	1.6	5.0%	8.1%	
Unexplned variance in 3rd contrast	=	1.5	4.6%	7.6%	
Unexplned variance in 4th contrast	=	1.3	4.1%	6.7%	
Unexplned variance in 5th contrast	=	1.2	3.7%	6.0%	
	• /•		1	•	

Figure 3 Standardized Residual variance (in Eigenvalue units)[A4]

Gradual Scale PCM analysis offers a unique verification process on the five grade categories of students' conceptual understanding (see Table 2). In Figure 4, it is illustrated that the average observation starts from logit -0.38 in category 1 (lack of knowledge) and increases up to logit +0.64 in category 5 (scientific knowledge). Such finding indicates that the grade category of students' conceptual understanding from 1 to 5 is considered as "very good". Moreover, PCM threshold is also employed to identify the grade's validity; the indicator highlights a transition that occurs in the students' decision making process from one grade to another (Linacre, 2012). The result of PCM-Andrich threshold analysis indicates a consistent

increase from the grade 1 to 5, implying that the grade category of conceptual understanding implemented as the evaluation scale to assess the students' competence is categorized as "very good".

5	SUMMAR	YO	F CATEGO	DRY S	STRUCTUR	E. Mo	del="R"				
ļ		ORY		RVED		AMPLE	INFIT O	UTFIT		CATEGORY	
						+		++-			
Ì	1	1	1743	20	38	34	.93	.99	NONE	(-1.39)	1 LOK (Lack of Knowledge)
İ	2	2	1173	14	.06	05	1.19	1.23	.21	46	2 AM (All-Misconception)
Ì	з	з	873	10	.21	.19	1.01	1.03	.37	.03	3 MFN (Misconception False Negativ
ĺ	4	4	1203	14	.33	.41	1.11	1.05	02	.49	4 MFN (Misconception False Positiv
	5	5	3548	42	.64	.63	1.00	1.04	56	(1.31)	5 SK (Scientific Knowledge)
C	BSERV	ED	AVERAGE	is n	nean of	measur	es in c	ategorv	. It is no	ot a param	neter estimate.

Figure 4. Validity of Grade Scale

Validity

The notion of validity revolves around the question: "does the test measure what it is supposed to measure?". That being said, the developed instrument is considered to have good construct validity if it is able to measure the students' conceptual understanding in explaining the concept of change of state of matter (Linacre, 2012, 2020; Sumintono, 2018).

Table 3

Item Statistics: Misfit Order

Itam	Maaguma	INI	FIT	OUT	FIT	PTMEA
nem	Measure	MNSQ	ZSTD	MNSQ	ZSTD	Corr.
1LG-1	86	1.65	4.6	1.62	3.4	.36
8SL-3	33	1.18	2.4	1.30	2.7	.44
11SG-1	29	1.17	2.3	1.29	2.8	.47
2LG-2	35	1.27	3.3	1.29	2.6	.44
5LG-5	.66	1.21	3.3	1.25	2.8	.47
12SG-2	21	1.10	1.5	1.14	1.5	.52
6SL-1	.04	1.03	.5	1.12	1.6	.40
19SL-4	37	1.08	1.1	.98	2	.56
7SL-2	.11	.95	9	1.08	1.2	.48
9SL-4	.04	1.01	.2	1.05	.7	.45
16LS-1	16	1.00	.1	1.02	.3	.51
4LG-4	.07	.97	5	1.01	.1	.48
17LS-2	40	1.00	.0	.92	7	.51
13SG-3	.30	.99	1	.96	5	.54
14SG-4	.04	.96	8	.88	-1.7	.61
15SG-5	.49	.91	-1.6	.86	-2.0	.55
3LG-3	.15	.85	-3.1	.91	-1.3	.48
20LS-5	.57	.80	-3.8	.90	-1.3	.50
10SL-5	.79	.73	-4.3	.78	-2.5	.48
18LS-3	29	.74	-4.0	.72	-3.2	.58

The first step is to ensure that all items match the Partial-Credit Model. Table 3 above displays the analysis result of statistic items. The study employs three criteria to measure any misfits or outliers between the students and the items (Linacre, 2012, 2020; Sumintono & Widhiarso, 2014): 1) the accepted Outfit Mean Square (MNSQ) value is between 0.5 < MNSQ < 1.5; 2) the accepted Outfit Z-Standard (ZSTD) value is between -2.0 < ZSTD < +2.0; 3) the accepted Point Measure Correlation (Pt Mean Corr) value is between 0.4 < Pt

Measure Corr. < 0.85. As illustrated in figure 4, it is detected that the item 1LG-1 does not meet two criteria (Outfit MNSQ and Outfit ZSTD), while the items 8SL-3, 11SG-1, 2LG-2, and 5LG-5 do not meet the Outfit ZSTD criteria. Moreover, the study does not find any items with negative results in the Point Measure Correlation criteria. This signifies that there is no single item that meets all the three criteria; thus, the measurement instrument possesses good item validity.

The second step is to measure the consistency between the item difficulty level and students' conceptual understanding. Figure 5 displays Wright Map to represent the item difficulty test level and students' conceptual understanding level. The graphical map is a result of empirical analysis on the answer response of the students in each item. According to the Wright Map result, all items in the measurement instrument has majorly encompassed the students' ability. The result indicates that the most difficult item is 10SL-5 (+0.79 logit), while the easiest item is 1LG-1 (-0.86 logit). However, no equivalent items were found in the understanding level smaller than -0.86 logit (-0.86 to -0.40 logit) as well as the level higher than +0.79 logit; therefore, further investigation is required.



Figure 5 Wright Map: Person-Map-Item (LG = liquid-gas, SL = solid-liquid, SG = solid-gas, LS = liquid-gas)

The research discovers several interesting cases regarding the difference between the items and students' conceptual understanding: Firstly, there are four items identified (LG, SL, SG and LS) that measure similar constructs within each level of conceptual understanding. Despite being in the same conceptual understanding level, the items' logit is completely different. For instance, four items were discovered in level 3, each with varying logit (8SL-3 (-0.33) < 18LS-3 (-0.29) < 3LG-3 (+0.15) < 13SG-3 (+0.30)). The numbers indicate that overall, students are more capable of explaining the concept of SL state change compared to LS, LG, and SG. This condition also occurs in the level 4, in which each item has varying logit (19LS-4 (-0.37) < 9SL-4 (+0.04) = 14SG-4(+0.04) < 4LG-4(+0.07)). Such a finding shows that the students find it easier to explain the correlation between the state of matter and the change process in LS compared to either SL, SG, or LG. Two sample cases above have illustrated that the students' conceptual understanding differs between the change process of LG (evaporation), SG (sublimation), SL (melting), and LS (freezing).

Moreover, it is found that the items in higher conceptual understanding levels tend to have lower logit than those at a lower level. As an instance, the logit of item 19SL-4 in level 4 (-0.37) is smaller than that of item 13SG-3 in level 3 (+0.30). This signifies that students find it harder to explain the item 13SG-3 compared to item 19SL-4. Thirdly, in the same concept of change of state (for example, LS), the logit of item 17LS-2 in level 2 (-0.40) is smaller than that of item 16LS-1 in level 1 (-0.16). As illustrated by the number, students find it easier to explain the SL concept in level 2 rather than to explain the concept's macroscopic fact in level 1. The findings above indicate that the students' conceptual understanding is not consistent with the item sequence. Moreover, the findings also suggest that the item difficulty level (LG, SL, SG, and LS; particularly SL (melting) and LS (freezing)) do not match the level in the construct map.

Measurement reliability

In Partial-Credit Model analysis, the indicator of reliability is observed from the quality of students' response patterns, the instrument, and the interaction between person-item. Within this study, item separation and person separation values are employed as the indicators. The separation index is also converted to Cronbach-equivalent value with an estimation of 0-1. The summary of measurement instrument statistics is displayed in Table 4.

	Student	Item
	(N=427)	(N=20)
Mean	0.26	0.00
Standard Error	0.02	0.09
Standard Deviation (SD)	0.48	0.41
Reliability	0.82	0.99
Infit mean-square	1.02	1.03
Outfit mean-square	1.05	1.05
Infit ZSTD	0.00	0.00
Outfit ZSTD	0.10	0.30
Point Raw Score to measure correlation	0.99	-0.99
Separation index (reliability)	2.10	9.54
Cronbach Alpha (KR-20): 0.84		
Data Points : 8540		
Chi-Square : 21173		
df : 8091 (p = 0.0000)		

Table 4.	Summary	of fit	statistics
----------	---------	--------	------------

From the table 4, it is generated that the total data points are 8540 with a Chi-square value of 21173 and the degree of freedom (df) of 8091 (p = 0.0000). These numbers indicate that the

measurement is deemed as "very good" and "significant". The column of students and item in the table 4 suggest whether or not the students and the item are considered fit. The average measure value of students is +0.26 logit (μ > 0.00), signifying that the students in overall are competent to explain the concept of change of state of matter. If the separation index value of students (+2.10 logit) is inputted into the person strata (H) formula, or H = [(4*separation) + 1]/3, thus, the generated H value = +3.13 (Linacre, 2012; Sumintono & Widhiarso, 2015). The person strata value (H) of 3 suggests that the students are classifiable into three groups of conceptual understanding (high, moderate, and low). On top of that, if the item's separation index value (+9.54) is processed by the same formula (H), the generated value is 13. Such a number shows that the items in the instrument are classifiable into 14 levels of difficulty. Moreover, the data illustrate that the items are deemed accurate and capable of measuring the students' competence in explaining the focused topic.

From the analysis result of students' answer pattern, the research generates Infit and Outfit MNSQ values of 1.02 and 1.05, respectively, with expectation value of 1.0. This clarifies that the students' answer pattern towards the instrument is categorized as "good" (Linacre, 2012; Sumintono & Widhiarso, 2014). In addition, the result generates Infit ZSTD and outfit ZSTD value of 0.0 and 0.10, respectively, with an expectation value of 0.0; the numbers depict that the overall students' answer pattern is in accordance with the model. Moreover, the overall reliability of students section is 0.82, categorized as "good". From the instrument item assessment, it is generated that the Infit and Outfit MNSQ values are 1.03 and 1.05, respectively, with the expectation value of 1.0, and the Infit and Outfit ZSTD values are 0.0 and 0.3, with the expectation value of 0.0. The numbers suggest that the overall instrument is deemed as "good", proven by the instrument reliability value of 0.99. The KR-20 (alpha Cronbach) value results in 0.84, thus signifying a good interaction between the students and the item. As acquired from the findings, the actual data in this study have met the Partial-Credit Model requirements, meaning that further analysis is considered as valid to conduct.

Level of Students' Learning Progress

The second problem of the research is: "How is the learning progress of the participants ranging from senior high school to fourth college year in explaining the focused topic?". To elaborate on that matter, the study employs data generated from the development process of 4TMC instrument to measure the students' conceptual understanding level.





(Senior high school students = A, first-year college students = B, second-year college students = C, third-year college students = D, fourth-year college students = E)

Figure 6 displays the average competence calculated in the form of logs based on the students' academic level, ranging from A to E. The figure shows an increasing trend in students' competence development based on their respective academic level (ABCDE). Moreover, it is discovered that the group E shows better learning progress compared to the other groups (D, C, B, and A). Despite that, the One-way ANOVA test indicates a difference among the students' competence based on the academic level, in which F_{count} (6, 0142442) > F_{table} (2,39308); df = 422; p <0.05. The research, therefore, conducted a post hoc Bonferroni test to identify which group that experience significant learning progress. As extracted from the statistical result, group A and B undergo significant learning progress, while group C, D, and E do not experience such significant advancement. This contradicts the common notion that the group CDE are college students with longer formal education experience compared to

group A or B. Such finding indicates that the group CDE find it hard to explain the concept of change of state of matter.

Comparison of average competence between groups ABCDE is conducted to map out the difference in the students' learning progress in each conceptual understanding level (displayed in Table 3). The students' competence is calculated based on four items in each level of conceptual understanding. As an example, in the level 1, the students' competence is measured by referring to the mean of item 1LG-1, 6SL-1, 11SG-1, and 16SL-1; the same also applies in the next levels. Based on Table 5, it is found that the students' competence in level 1 (0.77 logit, SD = 0.86) is higher than their competence in level 2 (0.69 logit, SD = 0.95); the same also applies in the next levels. The findings above indicate that the students' conceptual understanding has not developed optimally. On top of that, the item sequence in level 1 is easier to explain compared to that in level 2. The same condition also applies in the next levels. Students find it harder to explain concepts of change of state of matter as the learning progress level increases. Simply put, the students' learning progress level is different in each level of conceptual understanding.

Table 5

Measurement of students' average competence in each level of conceptual understanding

Conceptual Understanding		St	tudents' Educati	ion Level (Mean	, SD)	
Level	A (N=171)	B (N=83)	C (N=66)	D (N=55)	E (N=52)	ABDCE (N=427)
1	0.69 (0.86)	0.80 (0.71)	0.61 (0.91)	1.29 (0.95)	1.05 (0.90)	0.77 (0.86)
2	0.58 (1.04)	0.66 (0.75)	0.68 (0.86)	1.05 (1.00)	0.83 (0.97)	0.69 (0.95)
3	0.19 (0.95)	0.61 (1.00)	0.33 (1.13)	0.84 (0.92)	1.10 (1.24)	0.51 (1.10)
4	0.24 (1.00)	0.53 (0.68)	0.51 (1.12)	0.70 (0.86)	0.51 (0.71)	0.41 (0.57)
5	-1.16 (1.59)	-0.80 (1.46)	-0.86 (1.51)	-0.48 (0.85)	-0.58 (1.51)	-0.84 (1.41)

The difference in students' learning progress levels in each conceptual understanding level depicts that longer formal education experience does not necessarily guarantee that the student will have better learning progress in explaining the focused topic. For instance, Table 6 illustrates the comparison of item logit size in level 3 that is calculated based on the students' academic level.

Education	N	Item mean (logit) at level 3				
Level	14	13SG-3	3LG-3	18LS-3	8SL-3	
А	171	0.51	0.33	-0.22	-0.61	
В	83	0.55	0.40	-0.43	-0.51	
С	66	0.61	0.19	-0.15	-0.66	
D	55	0.33	0.20	-0.15	-0.46	
E	52	0.57	0.06	-0.30	-0.33	

Table 6 Average item logit in level 3

Table 7

Category of item 13SG-3 comprehension

		Conceptual Understanding Category - Item 13SG				
Grade	N					
		LOK	<mark>AM</mark>	<mark>MFN</mark>	MFP	<mark>SK</mark>
A	<mark>171</mark>	<mark>36</mark>	<mark>21</mark>	<mark>3</mark>	<mark>20</mark>	<mark>19</mark>
B	<mark>83</mark>	<mark>19</mark>	<mark>36</mark>	<mark>8</mark>	<mark>5</mark>	<mark>31</mark>
C	<mark>66</mark>	<mark>36</mark>	<mark>27</mark>	<mark>2</mark>	<mark>8</mark>	<mark>27</mark>
D	<mark>55</mark>	<mark>13</mark>	<mark>24</mark>	<mark>4</mark>	7	<mark>53</mark>
E E	<mark>52</mark>	<mark>23</mark>	<mark>12</mark>	<mark>6</mark>	<mark>12</mark>	<mark>48</mark>

Category: LOK = Lack of Knowledge, AM = All-Misconception, MFN = Misconception False Negative, MFP = Misconception False Positive, SK = Scientific Knowledge

How is the students' learning progress level in the same item? Table 7 displays the percentage data of students' competence in explaining item 13SG-3 based on five categories of conceptual understanding (LOK, AM, MFN, MFP, and SK). In the SK category, students in group D perform better among all groups (D (53%) > E (48%) > B (31%) > C (27%) > A (19%)). Simply put, more than half students in group D are capable of explaining the item 13SG-3 compared to students in other groups. Meanwhile, in LOK, students in group A and C show higher percentage among all groups (A (36%) = C (36%) > E (23%) > B (19%) > D (13%)). In other words, more than one-third of students in group A or C is incapable of explaining the item 13SG-3 compared to students in other groups B shows highest percentage among all groups (B (36%) > C (27%) > D (24%) > A (21%) > E (12%)); it signifies that more than one-fourth of students in group B are incapable of explaining item 13SG-3 compared to other

groups due to the misconception on the item. Such findings indicate that the high percentage in LOK and AM category is seen as one of the reasons why the students' competence is different in explaining the same item 13SG-3. To put it another way, the students' learning progress does not develop optimally in explaining item 13SG-3 due to lack of knowledge (LOK) or misconception (AM) on the item.



Figure 7(a)—__Probability Category Curve of item 13SG-3 of group A, and Figure 7(b) Probability Category Curve of item 13SG-3 of group D (Category: LOK = Lack of Knowledge, AM = All-Misconception, MFN = Misconception False Negative, MFP = Misconception False Positive, SK = Scientific Knowledge)

Figure 7 illustrates the comparison of the probability category curve (PCC) of students in group A and D in item 13SG-3. The five curve shapes are the visual representation of the distribution of five categories of students' conceptual understanding. From the curves, one can identify which groups that tend to show LOK and AM category traits. It is worth noting that the curve 2(a) and 2(b) tend to be different based on the MFP curve shape, while others are relatively similar. The MFP curve of students A has a higher probability compared to that of students D; simply put, a senior high school student tends to show stronger MFP category compared to a third-year college student. The notion is supported by the finding that senior high school students are relatively incapable of providing correct reason on item 13SG-3

compared to third-year college students. On the other hand, students with low ability in group D tend to show similar curve shape of LOK, AM, and MFN with group A. This implies that both groups' conceptual understanding in the item is relatively similar. In other words, the learning progress of group D, particularly in students with low ability, has not developed optimally despite the fact that that group D consists of third-year college students that progressed through three years of formal education experience in university.

Discussion

The result shows that: firstly, based on the logit size, the items are put in the following order: 13SG-3 > 3LG-3 > 18 SL-3 > 8SL-3. This is to say that it is harder for the students to explain the concept in item 13SG-3 compared to 3LG-3, 18SL-3, and 8SL-3. Secondly, the students' competence in each item is different and not in sequential order based on the education level (ABCDE). The finding leads to an assumption that all students in group E are supposed to perform better in explaining the item sequence in level 3 than those in group D, C, B, and A, since they progressed through longer education experience. However, the calculation result shows a different insight. In the item 13SG-3, students in group C are the most competent among all group (C (0.61) > E (0.57) > B (0.55) > A (0.51) > D (0.33)), while in the item 8SL-3, group E students are the most competent (E (-0.33) > D (-0.46) > B (-0.51) > A (-0.61) > C (-0.66)). Such a finding indicates that the students' competence is varied despite being at the same level. To put it another way, longer formal education tends to have an insignificant effect on the development of students' conceptual understanding.

This echoes previous findings that the learning progress is highly dependent on the students' learning process and experience (Duschl et al., 2011; Park et al., 2017; Wilson, 2009). Learning progress is defined as a sophisticated and systematic way of thinking, in which the students will undergo gradual progress when learning a topic for a long time interval. Students are able to ask questions, form hypotheses, design experiments to test hypotheses,

collect data, and draw conclusions (Sutiani et al., 2021). Such a systematic way of thinking is formed by the learning practices and education experience (Emden et al., 2018). Student's way of thinking is affected by learning experience, learning motivation, self regulation and self efficacy to explore understanding of how students go about learning (Haarala-Muhonen et al., 2016; Karagiannopoulou et al., 2020). On top of that, the research findings are in line with previous studies that highlighted that students have distinctive comprehension formed by their own experience (Chi et al., 2018; Emden et al., 2018; Hoe & Subramaniam, 2016; Jin, Mikeska, Hokayem, & Mavronikolas, 2019; Rogat et al., 2011; Testa et al., 2019). Such distinctive knowledge has not been explored by evaluation or intervention through learning roadmaps that are in accordance with remedial learning (Smith et al., 2006). In spite of that, it is considered essential to conduct a further analysis that focuses on the modification of conceptual understanding category and analysis variation that is able to define the characteristics of students' alternative conception. The development procedures, as explained in the methodology, has resulted in 4TMC instrument; however, instrument development is seen as an essential continuous process (Wilson, 2009, 2012).

Based on the research findings, the study identifies several important notes on the development of the 4TMC instrument. Firstly, further analysis of the characteristic of students' response behavior is necessary to conduct regarding the item clarity and the measured concept. The findings have implied that the percentage of LOK and AM understanding category is relatively dominant and tends to increase along with the level of conceptual understanding. Hence, the development of the concept level requires taking into consideration any potential term use that might confuse the students. A further study on the identification of commonly-understood terms or concepts is therefore essential. Secondly, a separate analysis is required to diagnose the factors contributing to the students' lack of knowledge and misconception. Regarding that, further analysis can be conducted by applying the analysis methods developed by previous studies (Caleon & Subramaniam, 2010; Hoe &

Subramaniam, 2016; Oon & Subramaniam, 2013). Thirdly, it is discovered that the concepts LG, SG, SL and LS were interpreted differently by the students. Despite being in the same conceptual understanding level, the items' difficulty level are completely different. Therefore, an evaluation on answer choices requires one to focus on the representation of understanding at the same level.

One of the features of the **Partial-Credit Model** is that the model facilitates one to identify any correlation between the construct map and the students' competence in ways that the students' competence can be analyzed by referring to the difference in item difficulty level. The 4TMC instrument indicates that there are students with very high ability as well as students with low ability in each group. Such a gap serves as the basis for qualitative interpretation to elaborate on the difference in students' competence. The insight is applicable in the learning process of chemistry subject. The instrument is expected to be beneficial for teachers in developing a formative test to identify the students' progress of conceptual understanding. On top of that, teachers are able to implement the instrument as a diagnostic instrument to evaluate students' conceptual understanding in providing feedback on their learning progress. Providing feedback also improves students' outcome and ability to undestand what they learn, increase students ability and creative thinking (Goulas & Megalokonomou, 2021; Redifer et al, 2021). Through this instrument teacher can give learning feedback to control students learning condition in learning environments both in theory and practice (Dijks et al., 2018; Latifi et al., 2021; Mahvelati, 2021). Further, the teachers will be able to develop instructional strategies that are specifically designed to tackle the students' difficulty in developing an epistemological explanation regarding the concept of change of state of matter. Through the development of this instructional strategies, teachers will be better able to focus on the goal orientation of learning achievement and motivate students to engange in learning activities (Lee & Keller, 2018; Guo & Leung, 2021; Lin et al., 2021)

Conclusions

The article elaborates on the development and validation procedures of the 4TMC instrument with Partial-Credit Model to evaluate the students' learning progress in explaining the concept of change of state of matter. In addition, the 4TMC instrument was tested on its effectiveness in providing reliable and valid information regarding students' conceptual understanding. The result revealed that the integration of the 4TMC test and Partial-Credit Model is effective and valid to be treated as the diagnostic instrument to measure students' learning progress. Moreover, it is discovered that students in group A, B, C, D, and E, particularly those with low ability, are hampered in developing an epistemological explanation of the concept. This blames the students' lack of certainty in their answer and reason; thus, assumed as having lack of knowledge or misconception. The low-ability students' curve shape of LOK and AM is consistent in the competence interval of less than 0.1 logit. On the other hand, the students' ability gets lower as the conceptual understanding level increases. Such finding indicates that the learning process and education experience provide a limited contribution for the students in developing a systematic way of thinking regarding the concept of change of state of matter. In spite of that, it is considered essential to conduct a further analysis that focuses on the modification of conceptual understanding category and analysis variation that is able to define the characteristics of students' alternative conception. The development procedures, as explained in the methodology, has resulted in 4TMC instrument; however, instrument development is seen as an essential continuous process.

Recommendations

The Based on the results of the study, there are several recommendations for researchers and teachers. For researchers, the findings of this research can be followed up to examine more in how students build their understanding gradually in explaining the concept of particles in substance form changes. The study can be conducted by developing tests that aim to evaluate

and diagnose the process of student knowledge formation and development while being able to identify at the level of education where the confusion of understanding occurs. The evaluation becomes more objective, not only reviewed from the student's point of ability but can be reviewed from the teacher's ability. The model of *PCM* multi-faced item response pattern approach becomes one of the important parts recommended for such objectives. In this way, students' ability to develop epistemological knowledge, and their ability to significantly actualize the knowledge gained can be measured well.

On the other hand, for teachers, the results of this study along with the stages of analysis approach used can be a reference in evaluating the progress of learners' learning, as well as determining alternative thinking frameworks of students in explaining the concept of substance change. The information serves as strategic feedback in formulating instructional strategies and preparing remedial learning, especially for students who have difficulty in developing epistemological explanations of substance changes.

Limitations

The limitations of the research are primarily related to the misrepresentation of student reasoning, which may arise in its efforts to connect phenomena and concepts measured in each item. In this context, the student may not excel to explain, because of his incapableness in using his heuristic reasoning. This instrument is not equipped with items that evaluate the heuristic abilities of the student in question. However, researchers decided to record this incompetence as a misconception or vague knowledge. For further research, it is recommended that the instrument be equipped with items that measure students' emotional and heuristic reasoning according to the conceptual framework to be evaluated.

Acknowledgments

The researchers would like to express their gratitude towards the Directorate of Research and Community Service, Ministry of Research and Technology of Republic of Indonesia, for the financial support through the University Basic Research Excellence Grant Program in the Research and Community Service Office of Universitas Negeri Gorontalo, 2020.

References [A5]

- Aktan, D. C. (2013). Investigation of students' intermediate conceptual understanding levels: The case of direct current electricity concepts. *European Journal of Physics*, *34*(1), 33–43. https://doi.org/10.1088/0143-0807/34/1/33
- Arslan, H. O., Cigdemoglu, C., & Moseley, C. (2012). A three-tier diagnostic test to assess pre-service teachers' misconceptions about global warming, greenhouse effect, ozone layer depletion, and acid rain. *International Journal of Science Education*, 34(11), 1667– 1686. https://doi.org/10.1080/09500693.2012.680618
- Bond, T. G., & Fox, C. M. (2007). *Applying The Rasch Model: Fundamental Measurent in the Human Sciences* (2nd Ed.). Routledge Taylor & Francis Group
- Caleon, I. S., & Subramaniam, R. (2010). Do students know what they know and what they don't know? Using a four-tier diagnostic test to assess the nature of students' alternative conceptions. *Research in Science Education*, 40(3), 313–337. https://doi.org/10.1007/s11165-009-9122-4
- Chandrasegaran, A. L., Treagust, D. F., & Mocerino, M. (2007). The development of two tier multiple-choice diagnostic instrument for evaluating secondary school students' ability to describe and explain chemical reactions using multiple levels of representation. *Chemistry Education Research and Practice*, 8(3), 293–307
- Chi, S., Wang, Z., Luo, M., Yang, Y., & Huang, M. (2018). Student progression on chemical symbol representation abilities at different grade levels (Grades 10–12) across gender. *Chemistry Education Research and Practice*, 19(4), 1055–1064. https://doi.org/10.1039/c8rp00010g
- Claesgens, J., Scalise, K., Wilson, M., & Stacy, A. (2009). Mapping student understanding in chemistry: The perspectives of chemists. *Science Education*, 93(1), 56–85. https://doi.org/10.1002/sce.20292.
- Djiks, M.A., Brummer, L., & Kostons, D. (2018). The anonymous reviewer: the relationship between perceived expertise and the perceptions of peer feedback in higher education. *Assessment & Evaluation in Higher Eduction*, 43(8), 1258-1271. https://doi.org/10.1080/02602938.2018.1447645 [TIDAK ADA DI MENDELEY SAYA]
- Duncan, R. G., & Hmelo-Silver, C. E. (2009). Learning progressions: Aligning curriculum, instruction, and assessment. *Journal of Research in Science Teaching*, 46(6), 606–609.

https://doi.org/10.1002/tea.20316

- Duschl, R., Maeng, S., & Sezen, A. (2011). Learning progressions and teaching sequences: A review and analysis. *Studies in Science Education*, 47(2), 123–182. https://doi.org/10.1080/03057267.2011.604476
- Emden, M., Weber, K., & Sumfleth, E. (2018). Evaluating a learning progression on "Transformation of Matter" on the lower secondary level. *Chemistry Education Research and Practice*, *19*(4), 1096–1116. https://doi.org/10.1039/c8rp00137e
- Goulas, S., & Megalokonomou, R. (2021). Knowing who you actually are: The effect of feedback on short-and longer-term outcomes. *Journal of Economic Behavior & Organization*, 183, 589-615. <u>https://doi.org/10.1016/j.jebo.2021.01.013</u>
- Guo, M., & Leung, F. K. S. (2020). Achievement goal orientations, learning strategies, and mathematics achievement: A comparison of Chinese Miao and Han students. *Psychology in the Schools*. <u>NEED VOLISS.PAGES.</u> https://doi.org10.1002/pits.22424_______
 [TIDAK ADA DI MENDELEY SAYA]
- Haarala-Muhonen, A., Ruohoniemi, M., Parpala, A., Komulainen, E., & Lindblom-Ylänne, S. (2016). How do the different study profiles of first-year students predict their study success, study progress and the completion of degrees? *Higher Education*, 74(6), 949– 962. https://doi.org/10.1007/s10734-016-0087-8
- Habiddin, & Page, E. M. (2019). Development and validation of a four-tier diagnostic instrument for chemical kinetics (FTDICK). *Indonesian Journal of Chemistry*, 19(3), 720–736. https://doi.org/10.22146/ijc.39218
- Hadenfeldt, J. C., Bernholt, S., Liu, X., Neumann, K., & Parchmann, I. (2013). Using ordered multiple-choice items to assess students' understanding of the structure and composition of matter. *Journal of Chemical Education*, 90(12), 1602–1608. https://doi.org/10.1021/ed3006192
- Hasan, S., Bagayoko, D., & Kelley, E. L. (1999). Misconceptions and the certainty of response index (CRI). *Physics Education*, *34*(5), 294–299. https://doi.org/10.1088/0031-9120/34/5/304
- Herrmann-Abell, C. F., & Deboer, G. E. (2016). Using rasch modeling and option probability curves to diagnose students' misconceptions. *American Eduacational Research Assossiation*, NEED VOLUSE 1–12

[PERBAIKAN / DIKOREKSI]

Herrmann-Abell, C. F., & Deboer, G. E. (2016). Using Rasch modeling and option probability curves to diagnose students' misconceptions. *Paper Presented at the 2016 American Eduacational Research Assossiation Annual Meeting Washington, DC April 8-12, 2016, 1–12.* https://files.eric.ed.gov/fulltext/ED572821.pdf

- Hoe, K. Y., & Subramaniam, R. (2016). On the prevalence of alternative conceptions on acidbase chemistry among secondary students: Insights from cognitive and confidence measures. *Chemistry Education Research and Practice*, 17(2), 263–282. https://doi.org/10.1039/c5rp00146c
- Jin, H., Mikeska, J. N., Hokayem, H., & Mavronikolas, E. (2019). Toward coherence in curriculum, instruction, and assessment: A review of learning progression literature. *Science Education*, 103(5), 1206–1234. https://doi.org/10.1002/sce.21525
- Karagiannopoulou, E., Milienos, F. S., & Rentzios, C. (2020). Grouping learning approaches and emotional factors to predict students' academic progress. *International Journal of School & Educational Psychology*, <u>NEED VOLISS. 1–</u> 18. https://doi.org.10.1080/2168363.2020.183241 [TIDAK ADA DI MENDELEY SAYA]
- Klassen, S. (2006). Contextual assessment in science education: Background, issues, and policy. *Science Education*, *90*(5), 820–851. https://doi.org/10.1002/sce.20150
- Latifi, S., Noroozi, O., & Talaee, E. (2021). Peer feedback or peer feedforward? Enhancing students' argumentative peer learning processes and outcomes. *British Journal of Educational Technology*, 52(2), 768-784. https://doi.org/10.1111/bjet.13054
- Lee, K., & Keller, J. M. (2021). Use of the ARCS model in education: A literature review. *Computers* & *Education*, *122*,<u>NEED</u> ISS. 54-62. https://doi.org/10.1016/j.compedu.2018.03.019 [TIDAK ADA DI MENDELEY SAYA]
- Lin, P. Y., Chai, C. S., Jong, M. S. Y., Dai, Y., Guo, Y., & Qin, J. (2021). Modeling the structural relationship among primary students' motivation to learn artificial intelligence. *Computers and Education: Artificial Intelligence*, 2, 100006NEED ISS.PAGES. [TIDAK ADA DI MENDELEY SAYA]
- Laliyo, Botutihe, & Panigoro. (2019). The development of two-tier instrument based on distractor to assess conceptual understanding level and student misconceptions in explaining redox reactions. *International Journal of Learning, Teaching and Educational Research*, 18(9), 216–237. https://doi.org/10.26803/ijlter.18.9.12
- Linacre, J. M. (2012). A user's guide to WINSTEPS® MINISTEP Rasch-model computer program: Program manual 3.75.0. <u>winsteps.com</u><u>NEED</u> <u>PUBLISHER.</u><u>https://doi.org/ISBN 0-941938-03-4</u> <u>DIKOREKSI</u>
- Linacre, J. M. (2020). A User's Guide to WINSTEPS® MINISTEP Rasch-Model Computer Programs Program Manual 4.5.1. <u>winsteps.comNEED</u>

PUBLISHER. https://doi.org/ISBN 0-941938-03-4 [PERBAIKAN / DIKOREKSI]

- Ling Lee, W., Chinna, K., & Sumintono, B. (2020). Psychometrics assessment of HeartQoL questionnaire: A Rasch analysis. *European Journal of Preventive Cardiology*, <u>NEED</u> 0(00)VOL.ISS.,1–6. https://doi.org/10.1177/2047487320902322____PERBAIKAN / DIKOREKSI
- Liu, X. (2012). Developing measurement instruments for science education research. In B. Fraser, K. G. Tobin, & C. J. McRobbie (Eds.), *Second international handbook of science education* (pp. 651–665). Springer Netherlands
- Lu, S., & Bi, H. (2016). Development of a measurement instrument to assess students' electrolyte conceptual understanding. *Chemistry Education Research and Practice*, 17(4), 1030–1040. https://doi.org/10.1039/c6rp00137h
- Mahvelati, E. H. (2021). Learners' perceptions and performance under peer versus teacher corrective feedback conditions. *Studies in Educational Evaluation*, 70. <u>NEED ISS.</u> <u>PAGES. https://doi.org/10.1016/j.stueduc.2021.100995</u> [TIDAK ADA DI MENDELEY SAYA]
- Morell, L., Collier, T., Black, P., & Wilson, M. (2017). A construct-modeling approach to develop a learning progression of how students understand the structure of matter. *Journal of Research in Science Teaching*, 54(8), 1024–1048. https://doi.org/10.1002/tea.21397
- Neumann, K., Viering, T., Boone, W. J., & Fischer, H. E. (2013). Towards a learning progression of energy. *Journal of Research in Science Teaching*, 50(2), 162–188. https://doi.org/10.1002/tea.21061
- Park, M., Liu, X., & Waight, N. (2017). Development of the connected chemistry as formative assessment pedagogy for high school chemistry teaching. *Journal of Chemical Education*, 94(3), 273–281. https://doi.org/10.1021/acs.jchemed.6b00299
- Peterson, R. F., Treagust, D. F., & Garnett, P. (1989). Development and application of a diagnostic instrument to evaluate grade-11 and -12 students' concepts of covalent bonding and structure following a course of instruction. *Journal of Research in Science Teaching*, *26*(4), 301–314. https://doi.org/10.1002/tea.3660260404
- Redifer, J. L., Bae, C. L., & Zhao, Q. (2021). Self-efficacy and performance feedback: Impacts on cognitive load during creative thinking. *Learning and Instruction*, 71. <u>NEED</u> <u>ISS. PAGES.</u> https://doi.org/10.1016/j.learninstruc.2020.101395 [TIDAK ADA DI MENDELEY SAYA]
- Rogat, A., Anderson, C., Foster, J., Goldberg, F., Hicks, J., Kanter, D., ... PROVIDE ALL

<u>AUTHORS UNTIL 20.</u> Wiser, M. (2011). Developing learning progression in support of the new science standards: A RAPID workshop series₃, <u>NEED VOL.</u> (4), 163-<u>NEED</u> <u>PAGES</u>. https://doi.org/10.12698/cpre.2011.lprapid

- Rogat, Aaron. (2011). Developing Learning Progressions in Support of the New Science Standards: A RAPID Workshop Series. CPRE Research Reports. Retrieved from http://repository.upenn.edu/cpre_researchreports/66 [HASIL PERBAIKAN / DIKOREKSI]
- Smith, C. L., Wiser, M., Anderson, C. W., & Krajcik, J. (2006). Implications of research on children's learning for standards and assessment: A proposed learning progression for matter and the atomic-molecular theory. *Measurement: Interdisciplinary Research & Perspective*, 4(1–2), 1–98. https://doi.org/10.1080/15366367.2006.9678570
- Sumintono, B., & Widhiarso, W. (2014). *Aplikasi model Rasch untuk penelitian ilmu-ilmu* sosial [Application of Rasch model in social science research]. Trim Komunikata. https://www.researchgate.net/publication/268688933%0AAplikasi

 Sutiani, A., Situmorang, M., & Silalahi, A. (2021). Implementation of an Inquiry Learning Model with Science Literacy to Improve Student Critical Thinking Skills. International Journal of Instruction, 14(2), 117-138. ARTICLE TITLE SHOULD BE LOWERCASE, JOURNAL TITLE AND VOULME SHOULD ITALIC. TIDAK ADA DI MENDELEY SAYA

- Testa, I., Capasso, G., Colantonio, A., Galano, S., Marzoli, I., Scotti di Uccio, U., ...<u>PROVUDE ALL AUTHORS UNTIL 20</u> Zappia, A. (2019). Development and validation of a university students' progression in learning quantum mechanics through exploratory factor analysis and Rasch analysis. *International Journal of Science Education*, 41(3), 388–417. https://doi.org/10.1080/09500693.2018.1556414
- Treagust, D. F. (1988). Development and use of diagnostic tests to evaluate students' misconceptions in science. *International Journal of Science Education*, 10(2), 159–169. https://doi.org/10.1080/0950069880100204
- Tyson, L., Treagust, D. F., & Bucat, R. B. (1999). The complexity of teaching and learning chemical equilibrium. *Journal of Chemical Education*, 76(2–4), 554–558. https://doi.org/10.1021/ed077p1560.1
- Wilson, M. (2005). *Constructing measures: an item response modeling approach*. Lawrence Erlbaum Associates, Inc. https://doi.org/10.4324/9781410611697
- Wilson, M. (2008). Cognitive diagnosis using item response models. <u>Journal of Psychology/</u> Zeitschrift Für Psychologie / Journal of Psychology, 216(2), 74–88. https://doi.org/10.1027/0044-3409.216.2.74
- Wilson, M. (2009). Measuring progressions: Assessment structures underlying a learning progression. *Journal of Research in Science Teaching*, 46(6), 716–730. https://doi.org/10.1002/tea.20318

Wilson, M. (2012). Responding to a challenge that learning progressions pose to measurement practice. In A. C. Alonzo & A. W. Gotwals (Eds.), *Learning progression in science* (pp. 317–344). <u>Sense PublishersNEED PUBLISHER</u>. https://doi.org/10.1007/978-94-6091-824-7 <u>HASIL PERBAIKAN /</u> DIKOREKSI



Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Congratulations! Your paper has been indexed by Scopus and ERIC (EU-JER ID#2011240749)

1 pesan

European Journal of Educational Research <editor@eu-jer.com> Balas Ke: European Journal of Educational Research <editor@eu-jer.com> Kepada: European Journal of Educational Research <syukrulhamdi@uny.ac.id> 18 Juli 2021 20.17

Dear Dr. Syukrul Hamdi,

As you remember, We have published your paper entitled "Implementation of Four-tier Instruments Based on the Rasch Model in Evaluating Students' Learning Progress" (Manuscript EU-JER ID#2011240749) in our April 2021 issue (Vol.10-Iss.2).

Congratulations! Your paper has been indexed by Scopus and ERIC.

Scopus link:

https://www.scopus.com/record/display.uri?eid=2-s2.0-85105034639&origin=inward&txGid= c8b889577d004f768d5deed50bb4fe2a&featureToggles=FEATURE_NEW_METRICS_SECTION:1

ERIC link:

https://eric.ed.gov/?q=Implementation+of+Four-Tier+Multiple-Choice+Instruments+Based+on+the+Partial+Credit+Model+in+Evaluating+Students%e2%80%99+Learning+Progress&id=EJ1294534

By the way, we need citations especially from different Scopus articles for the increasing our CiteScore (2.2). Could you announce your article to your colleagues in order to cite please?

We are looking forward to getting your and your colleagues contributions in the future.

Best regards,

Ahmet C. Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com



Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

URGENT REMINDER: The galley proof of your paper (ID#2007240601)

4 pesan

Editor - European Journal of Educational Research <editor@eu-jer.com> Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> 15 Oktober 2020 17.37

Dear Dr. Hamdi,

We will publish our new issue today.

Please get back urgently.

Best regards,

Ahmet C. Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

 Forwarded Message ------ Subject:REVISED EMAIL, PLEASE DON'T CONSIDER PREVIOUS ONE:The galley proof of your paper (ID#2007240601)
 Date:Tue, 13 Oct 2020 22:43:52 +0300
 From:Editor - European Journal of Educational Research <editor@eu-jer.com> To:Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

Dear Dr. Hamdi,

Please see the attached galley proof of your paper (ID#2007240601) (word file). Wrong parts have been highlighted in yellow in this finalized paper. Please highlight in green for your edited parts.

By the way,

- 1- The language of the paper should be edited by a native speaker as a proofreading lastly.
- 2- We couldn't use non English characters as our journal policy.
- 3- Please edit all references regarding with attached citation guide for APA 7 style. (Please see the citation guide page in our web site: https://www.eu-jer.com/citation-guide)

4- Please cite this article in order to improve your literature review:

Say, S., & Yildirim, F. S. (2020). Investigation of pre-service teachers' web 2.0 rapid content development selfefficacy belief levels and their views on web 2.0 tools. International Journal of Educational Methodology, 6(2), 345-354. https://doi.org/10.12973/ijem.6.2.345

We ask you to check it please. Please edit at word file and resend it to me please ASAP.

Best regards, Ahmet Savas Ph.D. Editor- European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

On 10/13/2020 7:39 PM, Syukrul Hamdi UNY wrote:

Dear, Editor of EUJER

I have recently revised (second round) the article according to the suggestions from reviewers 2 (R2612). I attached the result of the revised article. Thank you

Best regards Syukrul Hamdi

Pada tanggal Sel, 13 Okt 2020 pukul 15.26 Editor - European Journal of Educational Research <editor@eu-jer.com> menulis:

Dear Dr. Hamdi,

We have just received the feedback of the R2612. Please find the attached file as additional corrections. Please do these corrections and resend the finalized paper in 24 hours. Delete all old highlights and rehighlight for new edited parts.

We are looking forward to getting your finalized paper

Best regards,

Ahmet C. Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

On 13 Eki 2020 08:08, Syukrul Hamdi UNY wrote:

Dear, Editor of EU-JER

We are glad to receive the information regarding our article which is accepted for further publication. I inform you that I have already processed a payment through Bank Negara Indonesia for around USD 600. The staff of bank said that it will arrive at the Destination Bank in about 3-4 days. Thank you

Best regards,

Syukrul Hamdi

Pada tanggal Sen, 12 Okt 2020 pukul 21.51 Syukrul Hamdi UNY <<u>syukrulhamdi@uny.ac.id</u>> menulis:

Dear, Editor of EUJER

Thank you for the information. We are very happy to hear this information. We will make payment soon.

Best regards Syukrul Hamdi

Pada tanggal Sen, 12 Okt 2020 pukul 17.36 Editor - European Journal of Educational Research <editor@eu-jer.com> menulis:

Dear Dr. Syukrul Hamdi,

Congratulation! After a thorough double-blind review, I am pleased to inform you that your manuscript entitled "Development of Web-based Application for Teacher Candidate Competence Instruments: Preparing Professional Teachers in the IR 4.0 Era" (ID#2007240601) has been accepted. It is scheduled for publication in the Volume 9 Issue 4 of the "European Journal of Educational Research".

We kindly ask you to pay the article processing fee USD 500 and USD 100 transaction fee + tax of our bank (totally USD 600) via bank wire transfer. Kindly acknowledge invoice of this acceptance letter. Payment due date: October 13, 2020 in order to publish in our new issue.

Email Universitas Negeri Yogyakarta - URGENT REMINDER: The galley proof of your paper (ID#2007240601)

BANK WIRE TRANSFER INFORMATION : NAME OF BENEFICIARY: Ahmet Cezmi SAVAŞ ADDRESS OF BENEFICIARY: Degirmicem District Ozgurluk Str. No:32B, Zipcode:27090, Gaziantep, TURKEY PHONE OF BENEFICIARY: +90 (342) 909 61 90 CORRESPONDENT BANK CHARGER: REMITTER AMOUNT: USD 600 PAYMENT DETAIL: EU-JER Manuscript ID#2007240601 BANK NAME: QNB Finansbank BANK ADDRESS: Esentepe Mahallesi Büyükdere Caddesi Kristal Kule Binası No:215 Şişli - İstanbul BRANCH OF THE BANK: ENPARA BRANCH CODE: 3663 ACCOUNT NUMBER: 88177946 IBAN: TR66 0011 1000 0000 0088 1779 46 SWIFT CODE: FNNBTRISXXX

After payment, we will send the gallery proof of your paper. The galley proofs must be returned to us within 1 calendar day. Furthermore, you are responsible for any error in the published paper due to your oversight.

Thank you very much for submitting your article to the journal of "European Journal of Educational Research". We welcome your contributions in future.

Best regards.

Ahmet C. Savas Ph.D. Editor, European Journal of Educational Research http://www.eu-jer.com editor@eu-jer.com

Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya. (To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments) Universitas Negeri Yogyakarta

www.uny.ac.id

Untuk mendukung "Gerakan UNY Hijau", disarankan tidak mencetak email ini dan lampirannya. (To support the "Green UNY movement", it is recommended not to print the contents of this email and its attachments) Lipivorsitas Nogori Yogyakarta

Universitas Negeri Yogyakarta

www.uny.ac.id

Dear Dr.,

Thank you for your valuable feedback.

We are looking forward to getting your valuable contributions to our journal in the future.

Best regards,

Ahmet C. Savas, Ph.D. Editor, European Journal of Educational Research editor@eu-jer.com www.eu-jer.com

EU-JER_9_4_1751_HAMDI_PROOF.docx 1020K

Editor - European Journal of Educational Research <editor.eujer@gmail.com> Balas Ke: editor@eu-jer.com Kepada: Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id>

[Kutipan teks disembunyikan]

EU-JER_9_4_1751_HAMDI_PROOF.docx 1020K

Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> Kepada: European Journal of Educational Research <editor@eu-jer.com>

Dear Editor of Eu-Jer

I will send it soon. Thank you

Best Regards Syukrul Hamdi [Kutipan teks disembunyikan]

Syukrul Hamdi UNY <syukrulhamdi@uny.ac.id> Kepada: European Journal of Educational Research <editor@eu-jer.com>

Dear, Editor of EUJER

I have recently revised the article according to the suggestions. Thank you

Best regards Syukrul Hamdi [Kutipan teks disembunyikan]

EU-JER_9_4_1751_HAMDI_PROOF_Revised Version.docx 1020K 16 Oktober 2020 01.29

16 Oktober 2020 02.52

16 Oktober 2020 03.07

Copyright Transfer Agreement

European Journal of Educational Research [EU-JER] ("the Proprietor") will be pleased to publish your article ("the Work"), tentatively entitled

in the *EU-JER* ("the Journal") if the Work is accepted for publication. The undersigned authors transfer all copyright ownership in and relating to the Work, in all forms and media, to the Proprietor in the event that the Work is pub-lished. However, this agreement will be null and void if the Work is not published in the Journal.

The undersigned authors warrant that the Work is original, is not under consideration by another journal, and has not been previously published.

(This agreement must be signed by all authors. A photocopy of this form may be used if there are more than 10 authors.)

Author's name & signature	Date	Author's name & signature	Date
Author's name & signature	Date	Author's name & signature	Date
	Juie		
Author's name & signature	Date	Author's name & signature	Date
Author's name & signature	Date	Author's name & signature	Date
Author's name & signature	Date	Author's name & signature	Date

Copyright Transfer Agreement

European Journal of Educational Research [EU-JER] ("the Proprietor") will be pleased to publish your article ("the Work"), tentatively entitled

in the EU-JER ("the Journal") if the Work is accepted for publication. The undersigned authors transfer all copyright ownership in and relating to the Work, in all forms and media, to the Proprietor in the event that the Work is pub-lished. However, this agreement will be null and void if the Work is not published in the Journal.

The undersigned authors warrant that the Work is original, is not under consideration by another journal, and has not been previously published.

(This agreement must be signed by all authors. A photocopy of this form may be used if there are more than 10 authors.)

Author's name & signature Date Author's name & signature Date Mas 15090 Author's name & signature Date Author's name & signature Date ROMARIO AbdullaH Author's name & signature Date Author's name & signature Date Author's name & signature Date Date Author's name & signature Author's name & signature Date Date Author's name & signature



European Journal of Educational Research

Volume 10, Issue 2, 825 - 840.

ISSN: 2165-8714 http://www.eu-jer.com/

Implementation of Four-Tier Multiple-Choice Instruments <u>Based on the</u> <u>Partial Credit Model</u> in <u>Evaluating Students' Learning Progress</u>

Lukman Abdul Rauf Laliyo Universitas Negeri Gorontalo, INDONESIA **Syukrul Hamdi*** Universitas Negeri Yogyakarta, INDONESIA **Masrid Pikoli** Universitas Negeri Gorontalo, INDONESIA

Romario Abdullah Universitas Negeri Gorontalo, INDONESIA Citra Panigoro Universitas Negeri Gorontalo, INDONESIA

Received: May 5, 2020 • Revised: November 23, 2020 • Accepted: March 23, 2021

Abstract: One of the issues that hinder the students' learning progress is the inability to construct an epistemological explanation of a scientific phenomenon. Four-tier multiple-choice (hereinafter, 4TMC) instrument and Partial-Credit Model were employed to elaborate on the diagnosis process of the aforementioned problem. This study was to develop and implement the four-tier multiple-choice instrument with Partial-Credit Model to evaluate students' learning progress in explaining the conceptual change of state of matter. This research applied a development research referring to the test development model by Wilson. The data were obtained through development and validation techniques on 20 4TMC items tested to 427 students. On each item, the study applied diagnostic-summative assessment and certainty response index. The students' conceptual understanding level was categorized based on the combination of their answer choices; the measurement generated Partial-Credit Model for 1 parameter logistic (IPL) data. Analysis of differences was based on the student level class using Analysis of Variants (One-way ANOVA). This study resulted in 20 valid and reliable 4TMC instruments. The result revealed that the integration of 4TMC test and Partial-Credit Model was effective to be treated as the instrument to measure students' learning progress. One-way ANOVA test indicated the differences among the students' competence based on the academic level. On top of that, it was discovered that low-ability students showed slow progress due to the lack of knowledge as well as a misconception in explaining the Concept of Change of State of Matter. All in all, the research regarded that the diagnostic information was necessary for teachers in prospective development of learning strategies and evaluation of science learning.

Keywords: Learning progress, four-tier, change of state of matter, partial-credit model.

To cite this article: Laliyo, L.A.R., Hamdi, S., Pikoli, R., Abdullah, M., & Panigoro, C. (2021). Implementation of four-tier multiplechoice instruments based on the partial credit model in evaluating students' learning progress. *European Journal of Educational Research*, *10*(2), 825-840. https://doi.org/10.12973/eu-jer.10.2.825

Introduction

Central to the notion of science learning is the development of students' scientific understanding of basic concepts of sciences (Hadenfeldt et al., 2013), particularly, change of state of matter (Emden et al., 2018). Aside from the issue, several studies have also highlighted the students' inability to provide an epistemological explanation of basic concepts of sciences (Chi et al., 2018). Efforts to solve the issues, however, have shown little progress, as the students might have more complex perceptions regarding the alternative concept they understand (Morell et al., 2017).

Education practitioners have recommended the utilization of learning progress concept as the instructional method to provide guidance and direction and to adjust the curriculum with the learning process and assessment (Claesgens et al., 2009; Duncan & Hmelo-Silver, 2009; Rogat et al., 2011). Learning progress is defined as a sophisticated and systematic way of thinking. This method is applicable for a learning process, in which the students will undergo gradual progress when learning a topic in a long duration. Its effectiveness is highly dependent on the learning process and the students' learning experience (Duschl et al., 2011). The concept involves certain sets of gradual levels that represent conceptual understanding, ranging from low level up to comprehensive level.

The notion of learning progress is highly distinctive to each student and is dependent to one's learning experience (Rogat et al., 2011); therefore, there is no learning roadmap that is suitable for all kinds of students (Smith et al., 2006).

* Corresponding author:

© 2021 The Author(s). **Open Access** - This article is under the CC BY license (<u>https://creativecommons.org/licenses/by/4.0/</u>).

Syukrul Hamdi. Universitas Negeri Yogyakarta, Indonesia. 🖂 syukrulhamdi@uny.ac.id

Each student constructs one's understanding in a different way; moreover, the construction process is varied depending on the students' conceptual understanding level (Aktan, 2013). This is to say that each student undergoes a different rate of learning progress, understanding level, and knowledge construction. Simply put, the development of scientific comprehension among students is not linear (Neumann et al., 2013). Therefore, this study regards each level of students' conceptual understanding as a success in progressing for more advanced level of understanding (Hadenfeldt et al., 2013). A student who faces difficulty in a certain level of understanding will see a lack of progress to a more advanced level. This in turn hinders the student's ability to construct an epistemological explanation on the basic concepts of science. Within this context, the learning progress is treated as the method to evaluate students' conceptual understanding. The diagnostic information generated is reliable to be treated as a reference for the teachers in developing accurate and valid instructional components to guide the students to progress to the next level.

Among the diagnostic instruments that are considered applicable is the four-tier multiple-choice (4TMC) instrument. It is the development of two-tier multiple-choice test recommended by Treagust (1988) and Chandrasegaran et al., (2007). The use of two-tier instrument is familiar in identifying students' understanding in select topics such as electrochemistry (Lu & Bi, 2016), covalent bond (Peterson et al., 1989), and chemical equilibrium (Tyson et al., 1999). Despite its reputation in academia, the two-tier test has raised criticism due to its sole focus on the facts and negligence towards students' understanding (Klassen, 2006). Therefore, several experts propose the renewed version of the test by adding distractor answer choices to strengthen the diagnostic value of the items (Herrmann-Abell & DeBoer, 2016). In addition, some have highlighted the test's weakness in cases where students' tended to pick the answer choice and the reasoning randomly. This illustrates that the students were uncertain and possessed several misconceptions in the first tier question. In such cases, teachers faced difficulty in differentiating between guessed answers and misconceptions (Habiddin & Page, 2019; Hasan et al., 1999).

The criticism laid against the model has sparked the innovation of three-tier and four-tiers instruments. Both instruments feature two multi-level questions, also similar with two-tier test. In the three-tier test, however, the measurement of students' certainty level is conducted simultaneously in both first and second-tier questions; in the meantime, the measurement is conducted separately in the first two tiers (Caleon & Subramaniam, 2010). The value of students' certainty rate ranges from one (very uncertain) to five (very certain).

Three-tier test lacks validity in measuring the students' certainty rate regarding both the answer choice and the reasoning, whether or not the value of certainty rate refers only to the answer choice, to the reasoning, or both. Such weakness will in turn obstructs the evaluation and classification process of students' responses (Arslan et al., 2012). In the four-tier instrument, the measurement of certainty rate also involves the answer choice in the first tier and the reasoning in the third tier (Arslan et al., 2012). Regarding this feature, four-tier test is considered more accurate than the three-tier test. Students who pick wrong answer choices with high certainty indicate that they have a very high misconception on the measured item (Hoe & Subramaniam, 2016).

Four-tier instruments are used in studies discussing topics such as physics education (Caleon & Subramaniam, 2010), chemical thermodynamics (Sreenivasulu & Subramaniam, 2013), transition metal (Sreenivasulu & Subramaniam, 2013), acid-base reaction (Hoe & Subramaniam, 2016), and chemical kinetics (Habiddin & Page, 2019). However, it is worth noticing that studies on chemistry topic which employ four-tiers instruments tend to focus on describing alternative conception. To put it another way, the higher the certainty rate is, the stronger the students' alternative conception will be. Despite its potentials, the scholarly discussion has overlooked the implementation of a four-tier diagnostic instrument to measure students' learning progress. Therefore, further analysis is essential on the application of 4TMC test in several domains analyzes by Partial-Credit Model approach.

The use of Partial-Credit Model has been introduced since the 2000s in the science education research; it features the instrument that integrates diagnostic assessment and summative assessment (Liu, 2012). On top of that, the diagnostic assessment approach is introduced to conduct an in-depth analysis of the construction process of students' conceptual understanding (Claesgens et al., 2009; Hadenfeldt et al., 2013; Lu & Bi, 2016). This study employs 4MTC and Partial-Credit Model as a diagnostic tool to evaluate students' learning progress in explaining the change of state of matter, besides focusing on the Concept of Change of State of Matter, this research employs in-depth analysis using Item Response Theory, namely Partial Credit Model.

One of the features of the Partial-Credit Model is that the model facilitates one to identify any correlation between the construct map and the students' competence in ways that the students' competence can be analyzed by referring to the difference in item difficulty level. The 4TMC instrument indicates that there are students with very high ability as well as students with low ability in each group. Such a gap serves as the basis for qualitative interpretation to elaborate on the difference in students' competence. The insight is applicable in the learning process of chemistry subject. The instrument is expected to be beneficial for teachers in developing a formative test to identify the students' progress of conceptual understanding. On top of that, teachers are able to implement the instrument as a diagnostic instrument to evaluate students' conceptual understanding in providing feedback on their learning progress. Further, the teachers will be able to develop instructional strategies that are specifically designed to tackle the students' difficulty in developing an epistemological explanation regarding the concept of change of state of matter. The study focus revolves
around three research questions: 1) What is the quality of the developed 4TMC instrument based on the Partial-Credit Model?. 2). How is the effectiveness of 4TMC instrument to evaluate the students' learning progress in explaining concepts of change of state of matter. 3) How is the learning progress in students ranging from the senior high school level up to the senior (fourth) year of college in explaining the concepts?

Methodology

Development Model

This research used a development research referring to the test development model from Wilson. Wilson (2005, 2008) introduces four steps of measurement instrument development in figure 1.



Figure 1. Measurement instrument development

This recommendation is proven valid to be implemented in developing measurement instrument for different construct variables (Chi et al., 2018; Hadenfeldt et al., 2013; Laliyo et al., 2019; Lu & Bi, 2016; Wilson, 2009). The present study also included two questions related to certainty rate (Arslan et al., 2012; Habiddin & Page, 2019: Hasan et al., 1999). The obtained data were analyzed by Partial Credit Model (PCM) approach by WINSTEPS version 4.5.3 software.

Construct Map: Determining Level of Understanding

The first step was to develop the construct of measured variables. The study involved four concepts of change of state of matter: liquid-gas (LG), solid-liquid (SL), solid-gas (SG), and liquid-solid (LS). Gas-liquid (GL) and Gas-Solid (GS) materials were not included in this study as they are included in the basic level of knowledge. The change of a substance from gas to solid (GS) is known as freezing, while from gas to liquid (GL) is called condensing. These two types of changes in the form of substances are very easy to answer by students at a higher level since the materials have always been presented in textbooks, from high school to university students, on the topic of changes in the form of the substance. These concepts were implemented in a gradual manner through five levels of conceptual understanding (Table 1). Such method functions as the pathway of conceptual development that involves learning objectives from the lowest to the highest level of conceptual understanding (Duncan & Hmelo-Silver, 2009; Hadenfeldt et al., 2013; Rogat et al., 2011). In other words, the set of levels, as mentioned previously, was adjusted to the students' needs so as to develop their conceptual understanding; therefore, the levels, as illustrated in Table 1, was considered valid to illustrate the ideal conceptual development pathway (Neumann et al., 2013).

Cor	contuct Understanding Level	Change of State of Matter/Item					
Conceptual Understanding Level			SL	SG	LS		
5	Submicroscopic diagram of change of state of matter	5LG-5	10SL-5	15SG-5	20LS-5		
4	Correlation between state of matter and the process of change of state of matter	4LG-4	9SL-4	14SG-4	19LS-4		
3	Process of change of state of matter	3LG-3	8SL-3	13SG-3	18LS-3		
2	Concept of state of matter	2LG-2	7SL-2	12SG-2	17LS-2		
1	Factual phenomenon of state of matter	1LG-1	6SL-1	11SG-1	16LS-1		

Table 1. Level of conceptual understanding in explaining concept of change of state of matter

Description: (*LG* = *liquid-gas, SL* = *solid-liquid, SG* = *solid-gas, LS* = *liquid-gas*)

Item Design and Assessment Scheme

The second phase involved an item design. In the 4TMC instrument, all the items consisted of four-tier multiple-choices. To put it another way, each item contains four questions that combine between diagnostic-summative test (Hoe & Subramaniam, 2016; Lu & Bi, 2016; Treagust, 1988) with certainty response index (hereinafter, CRI) test (Arslan et al., 2012; Hasan et al., 1999). The first-tier questions (Q1) aimed to identify whether or not the students understand the content. Moreover, questions in the second tier (Q2) were employed to clarify the students' certainty regarding their answers in the Q1. Third-tier questions (Q3) functioned to diagnose the students' reasoning regarding their answers in the Q1. Further, questions in the second tier (Q4) were employed to clarify the students' certainty regarding their answers in the Q3. Q1 and Q3 questions in each item involved five answer choices; one among them was the correct answer, while three were the distractor, and another answer choice was open-ended answer choice. This open-ended option allows the students to decide the answer by themselves, should they find no correct answer as in accordance with their conceptual understanding. In the meantime, the Q2 and Q4 questions involved two close-ended answer choices; the first choice was for those who are uncertain of their answer, and the second choice was for the students who are very certain of their answer (Arslan et al., 2012). The distractor choices were employed in Q1 and Q3 questions to validate the diagnostic strength of the questions (Herrmann-Abell & DeBoer, 2016). Therefore, in the Q1 and Q3 tiers, the students would have only 0.20 or 20 percent probability of choosing the correct answer. The item Category of Grade of Students' Conceptual Understanding in Table 2.

Questions			Concontual	Understanding	Catagory	Dating	
Q1	Q2	Q3	Q4	conceptual	onderstanding	Category	Ratilig
Correct	Certain	Correct	Certain	CCCC	Scientific Knowledge	SK	5
Correct	Certain	Incorrect	Certain	CCIC	Misconception False Positi	ve MFP	4
Incorrect	Certain	Correct	Certain	ICCC	Misconception False Negative	MFN	3
Incorrect	Certain	Incorrect	Certain	ICIC	All-Misconception	AM	2
Correct	Certain	Correct	Uncertain	CCCU	Lack of Knowledge	LOK	1
Correct	Certain	Incorrect	Uncertain	CCIU	_		
Correct	Uncertain	Correct	Certain	CUCC			
Correct	Uncertain	Correct	Uncertain	CUCU			
Correct	Uncertain	Incorrect	Certain	CUIC			
Correct	Uncertain	Incorrect	Uncertain	CUIU			
Incorrect	Certain	Correct	Uncertain	ICCU			
Incorrect	Certain	Incorrect	Uncertain	ICIU			
Incorrect	Uncertain	Correct	Certain	IUCC			
Incorrect	Uncertain	Correct	Uncertain	IUCU			
Incorrect	Uncertain	Incorrect	Uncertain	IUIU			

Table 2 Category of grade of students' conceptual understanding*)

(*Hasan, Bagayoko and Kelley, 1999; Arslan, Cigdemoglu and Moseley, 2012; Habiddin and Page, 2019)

As an illustration, in the item 13SG-3, a student picks A in Q1, "very certain" in Q2, A in Q3, and "very certain" in Q4; the combination of the student's answers is ICIC. The result illustrates that the student's answer is incorrect in the Q1 and is very certain of one's error (Q2). Moreover, s/he also provides an incorrect answer in Q3 and is very certain of one's incorrect answer in Q3 (Q4). This indicates that in the item 13SG-3, the student is categorized to have all-misconception understanding (AM). In the Conceptual Understanding Category table, the category is included in fourth grade. Incorporation of the students' answer combinations in each item into the category and grade of students' understanding would result in specific data that are in accordance with the Partial-Credit Model.

Outcome Space and Data Collection

The third step involved the design of the outcome space of the correlation between items and construct maps (Bond & Fox, 2007; Wilson, 2009). The item validation was conducted independently by three expert validators to evaluate the extent of correlation between answer choices in Q1-Q3 in each item and the level of students' conceptual understanding. The validators were asked to clarify that the questions are easy to understand and the students' lack of linguistic competence would not hinder them from providing the right answer. The validators also required to ensure that the questions are in accordance with the syllabus, particularly with the students' conceptual understanding as based on the construct map. The questions in each item were also validated in several aspects, such as: ambiguity, time allocation, directiveness towards a particular answer, and subjective or emotional expression. Fleiss κ measure was employed to acquire information on the validators' approval. From the measure, it was generated that the κ value = 0.97, indicating that the three validators agreed that the 4TMC items were valid in correlating between the answer choices and the students' conceptual understanding.

The next step was to acquire data based on the measurement instrument. The instrument was tested to 427 students in Gorontalo, Indonesia using cluster random sampling technique. The students comprised 171 (40.05%) senior high

school students (or students A), 83 (19.44%) university freshmen majoring chemistry education (or students B), 66 (15.45%) second-year university students majoring chemistry education (or students C), 55 (12.88%) third-year university student majoring chemistry education (or students D), and 52 (12.18%) fourth-year university students majoring chemistry education (or students E). Based on gender, the female participants comprised 369 participants (86.41%), and the male counterparts consisted of 58 participants (13,58%). The participants were given no particular educational treatments and had stated their voluntary consent to participate in the research.

Partial-Credit Model Measurement and Data Analysis

The fourth step was to conduct the Partial-Credit Model measurement. This step was implemented to define the correlation between the score generated and the students' conceptual understanding level as elaborated within the construct map. The involvement of Partial-Credit Model measurement lay on the assumption that the item difficulty level is dependent on the students' answer, and that the students' understanding is dependent on the estimation of item difficulty (Linacre, 2012).

Partial credit model (PCM) was employed to evaluate the learning progress through structured questions; this took into account that the instrument items involved gradual and structured questions (Bond & Fox, 2007; Sumintono & Widhiarso, 2014; Wilson, 2009). The model was stated into the following formula: $\ln[P_nik/(1 - P_nik)] - B_n - D_ik$, in which P_{nik} refers to the probability of student n with B_n ability to pick correct response in the level k of item i; while D_{ik} refers to the difficulty level k of item i, or the threshold point for the test taker who scores k, not k -1. Analysis of differences based on class level of students using One-way ANOVA.

Results

The developed 4TMC instrument adapts the two-level instrument model by Treagust (1988), combined with the CRI theory by Arslan (2012). The function of CRI (certainty response index) is to ensure that students' choice of answers in Q1 and Q3 are the answers that they believe in. This is called diagnostic because it investigates the level of student error in stages, including the ability of students to understand and to use their understanding in explaining the reasons for their choice of answers. Thus, measurement is conducted both at the level of knowledge and reasoning.

The item design referred to the basic criteria to ensure that the students would be able to identify logical reason in Q3 as based on their answer in Q1; moreover, the item design also aimed to clarify the students' certainty of their answers through Q2 and Q4 questions. The 4TMC instrument also allows the students to state their certainty level of Q1 and Q3 answer choices separately. Students with correct understanding regarding the concept of change of state of matter (Q1) and its reasoning (Q3) will pick the "very certain" answer in the Q2 and Q4. If the students are uncertain of their answer regarding the content (Q1) but are certain of the reasoning (Q3), this suggests that the students are able to comprehend the concept/theory but unable to implement such concepts. This study views that it is beneficial to explore potential combinations of Q1/Q3 answer choices and Q2/Q4 certainty rate implementation to provide in-depth elaboration on students' understanding of certain concepts (Habiddin, 2019). The item design is illustrated in Figure 2.



Figure 2. 13SG-3 item design.

Each combination of students' answers in each item was categorized based on the assessment scheme. Every correct response in Q1 and Q3 is labeled with C (correct) code, and wrong answers were labeled with I (incorrect) code. Moreover, in the Q2 and Q4 tier, "very certain" answers are labeled with C (certain) code, and "uncertain" answers were labeled with U (uncertain) code. Therefore, combination of answer choice in each item was generated and written in sequence based on the questions in the Q1, Q2, Q3, and Q4. Such combination was treated as a reference in determining category and grade of students' conceptual understanding in each item, as shown in Table 2. The students' conceptual understanding category was adapted from findings reported by Hasan et al., (1999), Arslan et al., (2012), and Habiddin (2019).

Each combination of students' answer in each item was classified based on five categories of conceptual understanding. The first category is scientific knowledge (SK) with a grade of five; it illustrates that the students possess knowledge that is scientifically correct. The second category is misconception false positive (MFP) (with a grade of four), illustrating that the students have a correct claim of understanding, but they are unable to explain the claim. The third category is misconception false negative (MFN) (with a grade of three), illustrating that the students do not have correct claims of knowledge, but they are able to explain the claim. This category is considered negative because it is possible that the answers provided are guessed answer that is coincidentally correct. The fourth category is all-misconception (AM) with a grade of four; this category signifies that the students are very certain of their incorrect knowledge. Lastly, the fifth category is lack of knowledge (LOK) with a grade of five, signifying that the students lack knowledge in a particular item. Such categories were determined by the students' certainty level in Q2 or Q4. As an instance, one of the possible answer combinations is as follows: CCCU. This combination illustrates that the student is correct in Q1, Q2, and Q3, but is uncertain in Q4; the condition signifies that the student's understanding is ambiguous, hesitant, and is not based on appropriate scientific knowledge.

Effectiveness of Measurement Instruments

Unidimensionality is an essential indicator to evaluate the 4TMC instrument's ability to measure students' capability of explaining the concept of change of state of matter. This indicator is measured by Principal Component Analysis of the residuals to estimate the extent of variance to which the instrument is able to measure what it is supposed to measure (Sumintono & Widhiarso, 2014).

		Empirical	Modeled
Total raw variance in observations	=	32.7 100.0%	100.0%
Raw variance explained by measures	=	12.7 38.9%	39.2%
Raw variance explained by persons	=	3.8 11.7%	11.8%
Raw Variance explained by items	=	8.9 27.2%	27.5%
Raw unexplained variance (total)	=	20.0 61.1% 100.0	% 60.8%
Unexplned variance in 1st contrast	=	2.0 6.2% 10.2	%
Unexplned variance in 2nd contrast	=	1.6 5.0% 8.1	%
Unexplned variance in 3rd contrast	=	1.5 4.6% 7.6	%
Unexplned variance in 4th contrast	=	1.3 4.1% 6.7	%
Unexplned variance in 5th contrast	=	1.2 3.7% 6.0	%

Figure 3 Standardized residual variance (in eigenvalue units)

As displayed in the figure 3, the result of raw variance explained by measures of data is 38.9%, the number almost approaches the expectation value of 39.2%. The numbers indicate that the minimum unidimensionality requirements of 20% are achieved, and simultaneously, the limit of PCM unidimension is met (approaching 40%) (Linacre, 2012; Ling Lee et al., 2020). Moreover, the instrument's unexplained variance values are below 7% and considered as ideal (not exceeding 15%), signifying that the item independence rate in instrument falls into "good" category.

Gradual Scale PCM analysis offers a unique verification process on the five grade categories of students' conceptual understanding (see Table 2).

5	SUMM	ARY (OF CATEG	ORY S	STRUCTU	RE. Mo	del="R"					
ļ	CAT	EGOR	Y OBSE	RVED	OBSVD	SAMPLE	INFIT C	UTFIT	ANDRICH	CATEGOR	۲ļ	
	LAB	EL S	CORE COU	JNT %	AVRGE	EXPECT +	MNSQ	MNSQ	THRESHOLD	MEASUR	E -	
ĺ	1	1	1743	20	38	34	.93	.99	NONE	(-1.39	>j ı	LOK (Lack of Knowledge)
	2	2	1173	3 14	.06	05	1.19	1.23	.21	46	2	2 AM (All-Misconception)
	3	з	873	10	.21	.19	1.01	1.03	.37	.03	3	MFN (Misconception False Negativ
ĺ	4	4	1203	14	.33	.41	1.11	1.05	02	.49	14	MFN (Misconception False Positiv
ĺ	5	5	3548	3 42	.64	.63	1.00	1.04	56	j(1.31)İ ₅	SK (Scientific Knowledge)
C	DBSE	RVED	AVERAGE	is n	nean of	measur	es in c	ategory	/. It is n	ot a par	amet	er estimate.

Figure 4. Validity of Grade Scale

In Figure 4, it is illustrated that the average observation starts from logit -0.38 in category 1 (lack of knowledge) and increases up to logit +0.64 in category 5 (scientific knowledge). Such finding indicates that the grade category of students' conceptual understanding from 1 to 5 is considered as "very good". Moreover, PCM threshold is also employed to identify the grade's validity; the indicator highlights a transition that occurs in the students' decision making process from one grade to another (Linacre, 2012). The result of PCM-Andrich threshold analysis indicates a consistent increase from the grade 1 to 5, implying that the grade category of conceptual understanding implemented as the evaluation scale to assess the students' competence is categorized as "very good".

Validity

The notion of validity revolves around the question: "does the test measure what it is supposed to measure?". That being said, the developed instrument is considered to have good construct validity if it is able to measure the students' conceptual understanding in explaining the concept of change of state of matter (Linacre, 2012, 2020; Sumintono & Widhiarso, 2014).

		Tuble 5.1	cem statisti	es. 1.115/10 010		
Itom	Моодимо	INFIT		OUTFIT		РТМЕА
Item	Measure	MNSQ	ZSTD	MNSQ	ZSTD	Corr.
1LG-1	86	1.65	4.6	1.62	3.4	.36
8SL-3	33	1.18	2.4	1.30	2.7	.44
11SG-1	29	1.17	2.3	1.29	2.8	.47
2LG-2	35	1.27	3.3	1.29	2.6	.44
5LG-5	.66	1.21	3.3	1.25	2.8	.47
12SG-2	21	1.10	1.5	1.14	1.5	.52
6SL-1	.04	1.03	.5	1.12	1.6	.40
19SL-4	37	1.08	1.1	.98	2	.56
7SL-2	.11	.95	9	1.08	1.2	.48
9SL-4	.04	1.01	.2	1.05	.7	.45
16LS-1	16	1.00	.1	1.02	.3	.51
4LG-4	.07	.97	5	1.01	.1	.48
17LS-2	40	1.00	.0	.92	7	.51
13SG-3	.30	.99	1	.96	5	.54
14SG-4	.04	.96	8	.88	-1.7	.61

Table 3. Item statistics:	Misfit	order
---------------------------	--------	-------

Itom	Moacuro	INFIT	INFIT			PTMEA
nem	Measure	MNSQ	ZSTD	MNSQ	ZSTD	Corr.
15SG-5	.49	.91	-1.6	.86	-2.0	.55
3LG-3	.15	.85	-3.1	.91	-1.3	.48
20LS-5	.57	.80	-3.8	.90	-1.3	.50
10SL-5	.79	.73	-4.3	.78	-2.5	.48
18LS-3	29	.74	-4.0	.72	-3.2	.58

Table 3. Continued

The first step is to ensure that all items match the Partial-Credit Model. Table 3 displays the analysis result of statistic items. The study employs three criteria to measure any misfits or outliers between the students and the items (Linacre, 2012, 2020; Sumintono & Widhiarso, 2014): 1) the accepted Outfit Mean Square (MNSQ) value is between 0.5 < MNSQ < 1.5; 2) the accepted Outfit Z-Standard (ZSTD) value is between -2.0 < ZSTD < +2.0; 3) the accepted Point Measure Correlation (Pt Mean Corr) value is between 0.4 < Pt Measure Corr. < 0.85. As illustrated in figure 4, it is detected that the item 1LG-1 does not meet two criteria (Outfit MNSQ and Outfit ZSTD), while the items 8SL-3, 11SG-1, 2LG-2, and 5LG-5 do not meet the Outfit ZSTD criteria. Moreover, the study does not find any items with negative results in the Point Measure Correlation criteria. This signifies that there is no single item that meets all the three criteria; thus, the measurement instrument possesses good item validity.

The second step is to measure the consistency between the item difficulty level and students' conceptual understanding.



Figure 5 Wright Map: Person-Map-Item (LG = liquid-gas, SL = solid-liquid, SG = solid-gas, LS = liquid-gas)

Figure 5 displays Wright Map to represent the item difficulty test level and students' conceptual understanding level. The graphical map is a result of empirical analysis on the answer response of the students in each item. According to the Wright Map result, all items in the measurement instrument has majorly encompassed the students' ability. The

result indicates that the most difficult item is 10SL-5 (+0.79 logit), while the easiest item is 1LG-1 (-0.86 logit). However, no equivalent items were found in the understanding level smaller than -0.86 logit (-0.86 to -0.40 logit) as well as the level higher than +0.79 logit; therefore, further investigation is required.

The research discovers several interesting cases regarding the difference between the items and students' conceptual understanding: Firstly, there are four items identified (LG, SL, SG and LS) that measure similar constructs within each level of conceptual understanding. Despite being in the same conceptual understanding level, the items' logit is completely different. For instance, four items were discovered in level 3, each with varying logit (8SL-3 (-0.33) < 18LS-3 (-0.29) < 3LG-3 (+0.15) < 13SG-3 (+0.30)). The numbers indicate that overall, students are more capable of explaining the concept of SL state change compared to LS, LG, and SG. This condition also occurs in the level 4, in which each item has varying logit (19LS-4 (-0.37) < 9SL-4 (+0.04) = 14SG-4(+0.04) < 4LG-4 (+0.07)). Such a finding shows that the students find it easier to explain the correlation between the state of matter and the change process in LS compared to either SL, SG, or LG. Two sample cases above have illustrated that the students' conceptual understanding differs between the change process of LG (evaporation), SG (sublimation), SL (melting), and LS (freezing).

Moreover, it is found that the items in higher conceptual understanding levels tend to have lower logit than those at a lower level. As an instance, the logit of item 19SL-4 in level 4 (-0.37) is smaller than that of item 13SG-3 in level 3 (+0.30). This signifies that students find it harder to explain the item 13SG-3 compared to item 19SL-4. Thirdly, in the same concept of change of state (for example, LS), the logit of item 17LS-2 in level 2 (-0.40) is smaller than that of item 16LS-1 in level 1 (-0.16). As illustrated by the number, students find it easier to explain the SL concept in level 2 rather than to explain the concept's macroscopic fact in level 1. The findings above indicate that the students' conceptual understanding is not consistent with the item sequence. Moreover, the findings also suggest that the item difficulty level (LG, SL, SG, and LS; particularly SL (melting) and LS (freezing)) do not match the level in the construct map.

Measurement reliability

Chi-Square : 21173 df : 8091 (p = 0.0000)

In Partial-Credit Model analysis, the indicator of reliability is observed from the quality of students' response patterns, the instrument, and the interaction between person-item. Within this study, item separation and person separation values are employed as the indicators. The separation index is also converted to Cronbach-equivalent value with an estimation of 0-1. The summary of measurement instrument statistics is displayed in Table 4.

	Student (N=427)	Item (N=20)
Mean	0.26	0.00
Standard Error	0.02	0.09
Standard Deviation (SD)	0.48	0.41
Reliability	0.82	0.99
Infit mean-square	1.02	1.03
Outfit mean-square	1.05	1.05
Infit ZSTD	0.00	0.00
Outfit ZSTD	0.10	0.30
Point Raw Score to measure correlation	0.99	-0.99
Separation index (reliability)	2.10	9.54
Cronbach's alpha (KR-20): 0.84		
Data Points : 8540		

Table 4.	Summary	of fit	statistics
----------	---------	--------	------------

From the Table 4, it is generated that the total data points are 8540 with a Chi-square value of 21173 and the degree of freedom (df) of 8091 (p = 0.0000). These numbers indicate that the measurement is deemed as "very good" and "significant". The column of students and item in the table 4 suggest whether or not the students and the item are considered fit. The average measure value of students is +0.26 logit ($\mu > 0.00$), signifying that the students in overall are competent to explain the concept of change of state of matter. If the separation index value of students (+2.10 logit) is inputted into the person strata (H) formula, or H = [(4*separation) + 1]/3, thus, the generated H value = +3.13 (Linacre, 2012; Sumintono & Widhiarso, 2014). The person strata value (H) of 3 suggests that the students are classifiable into three groups of conceptual understanding (high, moderate, and low). On top of that, if the item's separation index value (+9.54) is processed by the same formula (H), the generated value is 13. Such a number shows that the items in the instrument are classifiable into 14 levels of difficulty. Moreover, the data illustrate that the items are deemed accurate and capable of measuring the students' competence in explaining the focused topic.

From the analysis result of students' answer pattern, the research generates infit and outfit MNSQ values of 1.02 and 1.05, respectively, with expectation value of 1.0. This clarifies that the students' answer pattern towards the instrument

is categorized as "good" (Linacre, 2012; Sumintono & Widhiarso, 2014). In addition, the result generates Infit ZSTD and outfit ZSTD value of 0.0 and 0.10, respectively, with an expectation value of 0.0; the numbers depict that the overall students' answer pattern is in accordance with the model. Moreover, the overall reliability of students section is 0.82, categorized as "good". From the instrument item assessment, it is generated that the infit and outfit MNSQ values are 1.03 and 1.05, respectively, with the expectation value of 1.0, and the infit and outfit ZSTD values are 0.0 and 0.3, with the expectation value of 0.0. The numbers suggest that the overall instrument is deemed as "good", proven by the instrument reliability value of 0.99. The KR-20 (Cronbach's alpha) value results in 0.84, thus signifying a good interaction between the students and the item. As acquired from the findings, the actual data in this study have met the Partial-Credit Model requirements, meaning that further analysis is considered as valid to conduct.

Level of Students' Learning Progress

The second problem of the research is: "How is the learning progress of the participants ranging from senior high school to fourth college year in explaining the focused topic?". To elaborate on that matter, the study employs data generated from the development process of 4TMC instrument to measure the students' conceptual understanding level.





(Senior high school students = A, first-year college students = B, second-year college students = C, third-year college students = D, fourth-year college students = E)

Figure 6 displays the average competence calculated in the form of logs based on the students' academic level, ranging from A to E. The figure shows an increasing trend in students' competence development based on their respective academic level (ABCDE). Moreover, it is discovered that the group E shows better learning progress compared to the other groups (D, C, B, and A). Despite that, the One-way ANOVA test indicates a difference among the students' competence based on the academic level, in which F_{count} (6, 0142442) > F_{table} (2,39308); df = 422; p <0.05. The research, therefore, conducted a post hoc Bonferroni test to identify which group that experience significant learning progress. As extracted from the statistical result, group A and B undergo significant learning progress, while group C, D, and E do not experience such significant advancement. This contradicts the common notion that the group CDE are college students with longer formal education experience compared to group A or B. Such finding indicates that the group CDE find it hard to explain the concept of change of state of matter.

Comparison of average competence between groups ABCDE is conducted to map out the difference in the students' learning progress in each conceptual understanding level (displayed in Table 3). The students' competence is calculated based on four items in each level of conceptual understanding. As an example, in the level 1, the students' competence is measured by referring to the mean of item 1LG-1, 6SL-1, 11SG-1, and 16SL-1; the same also applies in the next levels.

Table 5 Measurement of students' average competence in each level of conceptual understanding

Conceptual						
Level	A (N=171)	B (N=83)	C (N=66)	D (N=55)	E (N=52)	ABDCE (N=427)
1	0.69 (0.86)	0.80 (0.71)	0.61 (0.91)	1.29 (0.95)	1.05 (0.90)	0.77 (0.86)
2	0.58 (1.04)	0.66 (0.75)	0.68 (0.86)	1.05 (1.00)	0.83 (0.97)	0.69 (0.95)
3	0.19 (0.95)	0.61 (1.00)	0.33 (1.13)	0.84 (0.92)	1.10 (1.24)	0.51 (1.10)
4	0.24 (1.00)	0.53 (0.68)	0.51 (1.12)	0.70 (0.86)	0.51 (0.71)	0.41 (0.57)
5	-1.16 (1.59)	-0.80 (1.46)	-0.86 (1.51)	-0.48 (0.85)	-0.58 (1.51)	-0.84 (1.41)

Based on Table 5, it is found that the students' competence in level 1 (0.77 logit, SD = 0.86) is higher than their competence in level 2 (0.69 logit, SD = 0.95); the same also applies in the next levels. The findings above indicate that the students' conceptual understanding has not developed optimally. On top of that, the item sequence in level 1 is easier to explain compared to that in level 2. The same condition also applies in the next levels. Students find it harder to explain concepts of change of state of matter as the learning progress level increases. Simply put, the students' learning progress level is different in each level of conceptual understanding.

The difference in students' learning progress levels in each conceptual understanding level depicts that longer formal education experience does not necessarily guarantee that the student will have better learning progress in explaining the focused topic. For instance, Table 6 illustrates the comparison of item logit size in level 3 that is calculated based on the students' academic level.

Education	N	Item me	Item mean (logit) at level 3					
Level	IN	13SG-3	3LG-3	18LS-3	8SL-3			
А	171	0.51	0.33	-0.22	-0.61	-		
В	83	0.55	0.40	-0.43	-0.51			
С	66	0.61	0.19	-0.15	-0.66			
D	55	0.33	0.20	-0.15	-0.46			
E	52	0.57	0.06	-0.30	-0.33			

Table 6. Average item logit in level 3

Crada	N	Concept	Conceptual Understanding Category - Item 13SG-3 (%)					
Graue		LOK	AM	MFN	MFP	SK		
А	171	36	21	3	20	19		
В	83	19	36	8	5	31		
С	66	36	27	2	8	27		
D	55	13	24	4	7	53		
Е	52	23	12	6	12	48		

Table 7	Catagor	u of itom	1250 2	comprehension	•
i ubie 7.	calegor	v oj ilem	1336-3	comprehension	L

Category: LOK = Lack of Knowledge, AM = All-Misconception, MFN = Misconception False

Negative, MFP = Misconception False Positive, SK = Scientific Knowledge

How is the students' learning progress level in the same item? Table 7 displays the percentage data of students' competence in explaining item 13SG-3 based on five categories of conceptual understanding (LOK, AM, MFN, MFP, and SK). In the SK category, students in group D perform better among all groups (D (53%) > E (48%) > B (31%) > C (27%) >A (19%)). Simply put, more than half students in group D are capable of explaining the item 13SG-3 compared to students in other groups. Meanwhile, in LOK, students in group A and C show higher percentage among all groups (A (36%) = C (36%) > E (23%) > B (19%) > D (13%)). In other words, more than one-third of students in group A or C is incapable of explaining the item 13SG-3 compared to students in other groups due to the limited knowledge on the item. Moreover, in AM, group B shows highest percentage among all groups (B (36%) > C (27%) > D (24%) > A (21%) > E (12%)); it signifies that more than one-fourth of students in group B are incapable of explaining item 13SG-3 compared to other groups due to the misconception on the item. Such findings indicate that the high percentage in LOK and AM category is seen as one of the reasons why the students' competence is different in explaining item 13SG-3 due to lack of knowledge (LOK) or misconception (AM) on the item.



Figure 7. Probability Category Curve of item 13SG-3 of group A, and Probability Category Curve of item 13SG-3 of group D (Category: LOK = Lack of Knowledge, AM = All-Misconception, MFN = Misconception False Negative, MFP = Misconception False Positive, SK = Scientific Knowledge)

Figure 7 illustrates the comparison of the probability category curve (PCC) of students in group A and D in item 13SG-3. The five curve shapes are the visual representation of the distribution of five categories of students' conceptual understanding. From the curves, one can identify which groups that tend to show LOK and AM category traits. It is worth noting that the curve 2(a) and 2(b) tend to be different based on the MFP curve shape, while others are relatively similar. The MFP curve of students A has a higher probability compared to that of students D; simply put, a senior high school student tends to show stronger MFP category compared to a third-year college student. The notion is supported by the finding that senior high school students are relatively incapable of providing correct reason on item 13SG-3 compared to third-year college students. On the other hand, students with low ability in group D tend to show similar curve shape of LOK, AM, and MFN with group A. This implies that both groups' conceptual understanding in the item is relatively similar. In other words, the learning progress of group D, particularly in students with low ability, has not developed optimally despite the fact that that group D consists of third-year college students that progressed through three years of formal education experience in university.

Discussion

The result shows that: firstly, based on the logit size, the items are put in the following order: 13SG-3 > 3LG-3 > 18 SL-3 > 8SL-3. This is to say that it is harder for the students to explain the concept in item 13SG-3 compared to 3LG-3, 18SL-3, and 8SL-3. Secondly, the students' competence in each item is different and not in sequential order based on the education level (ABCDE). The finding leads to an assumption that all students in group E are supposed to perform better in explaining the item sequence in level 3 than those in group D, C, B, and A, since they progressed through longer education experience. However, the calculation result shows a different insight. In the item 13SG-3, students in group C are the most competent among all group (C (0.61) > E (0.57) > B (0.55) > A (0.51) > D (0.33)), while in the item 8SL-3, group E students are the most competent (E (-0.33) > D (-0.46) > B (-0.51) > A (-0.61) > C (-0.66)). Such a finding indicates that the students' competence is varied despite being at the same level. To put it another way, longer formal education tends to have an insignificant effect on the development of students' conceptual understanding.

This echoes previous findings that the learning progress is highly dependent on the students' learning process and experience (Duschl et al., 2011; Park et al., 2017; Wilson, 2009). Learning progress is defined as a sophisticated and systematic way of thinking, in which the students will undergo gradual progress when learning a topic for a long time interval. Students are able to ask questions, form hypotheses, design experiments to test hypotheses, collect data, and draw conclusions (Sutiani et al., 2021). Such a systematic way of thinking is formed by the learning practices and education experience (Emden et al., 2018). Student's way of thinking is affected by learning experience, learning motivation, self regulation and self efficacy to explore understanding of how students go about learning (Haarala-Muhonen et al., 2016; Karagiannopoulou et al., 2020). On top of that, the research findings are in line with previous studies that highlighted that students have distinctive comprehension formed by their own experience (Chi et al., 2018; Emden et al., 2018; Hoe & Subramaniam, 2016; Jin et al., 2019; Rogat et al., 2011; Testa et al., 2019). Such distinctive knowledge has not been explored by evaluation or intervention through learning roadmaps that are in accordance with remedial learning (Smith et al., 2006). In spite of that, it is considered essential to conduct a further analysis that focuses on the modification of conceptual understanding category and analysis variation that is able to define the characteristics of students' alternative conception. The development procedures, as explained in the methodology, has

resulted in 4TMC instrument; however, instrument development is seen as an essential continuous process (Wilson, 2009, 2012).

Based on the research findings, the study identifies several important notes on the development of the 4TMC instrument. Firstly, further analysis of the characteristic of students' response behavior is necessary to conduct regarding the item clarity and the measured concept. The findings have implied that the percentage of LOK and AM understanding category is relatively dominant and tends to increase along with the level of conceptual understanding. Hence, the development of the concept level requires taking into consideration any potential term use that might confuse the students. A further study on the identification of commonly-understood terms or concepts is therefore essential. Secondly, a separate analysis is required to diagnose the factors contributing to the students' lack of knowledge and misconception. Regarding that, further analysis can be conducted by applying the analysis methods developed by previous studies (Caleon & Subramaniam, 2010; Hoe & Subramaniam, 2016). Thirdly, it is discovered that the concepts LG, SG, SL and LS were interpreted differently by the students. Despite being in the same conceptual understanding level, the items' difficulty levels are completely different. Therefore, an evaluation on answer choices requires one to focus on the representation of understanding at the same level.

One of the features of the Partial-Credit Model is that the model facilitates one to identify any correlation between the construct map and the students' competence in ways that the students' competence can be analyzed by referring to the difference in item difficulty level. The 4TMC instrument indicates that there are students with very high ability as well as students with low ability in each group. Such a gap serves as the basis for qualitative interpretation to elaborate on the difference in students' competence. The insight is applicable in the learning process of chemistry subject. The instrument is expected to be beneficial for teachers in developing a formative test to identify the students' progress of conceptual understanding. On top of that, teachers are able to implement the instrument as a diagnostic instrument to evaluate students' conceptual understanding in providing feedback on their learning progress. Providing feedback also improves students' outcome and ability to understand what they learn, increase students ability and creative thinking (Goulas & Megalokonomou, 2021). Through this instrument teacher can give learning feedback to control students learning condition in learning environments both in theory and practice (Dijks et al., 2018; Latifi et al., 2021). Further, the teachers will be able to develop instructional strategies that are specifically designed to tackle the students' difficulty in developing an epistemological explanation regarding the concept of change of state of matter. Through the development of these instructional strategies, teachers will be better able to focus on the goal orientation of learning achievement and motivate students to engage in learning activities (Lee & Keller, 2021; Guo & Leung, 2021; Lin et al., 2021)

Conclusions

The article elaborates on the development and validation procedures of the 4TMC instrument with Partial-Credit Model to evaluate the students' learning progress in explaining the concept of change of state of matter. In addition, the 4TMC instrument was tested on its effectiveness in providing reliable and valid information regarding students' conceptual understanding.

The result revealed that the integration of the 4TMC test and Partial-Credit Model is effective and valid to be treated as the diagnostic instrument to measure students' learning progress. Moreover, it is discovered that students in group A, B, C, D, and E, particularly those with low ability, are hampered in developing an epistemological explanation of the concept. This blames the students' lack of certainty in their answer and reason; thus, assumed as having lack of knowledge or misconception. The low-ability students' curve shape of LOK and AM is consistent in the competence interval of less than 0.1 logit. On the other hand, the students' ability gets lower as the conceptual understanding level increases. Such finding indicates that the learning process and education experience provide a limited contribution for the students in developing a systematic way of thinking regarding the concept of change of state of matter. In spite of that, it is considered essential to conduct a further analysis that focuses on the modification of conceptual understanding category and analysis variation that is able to define the characteristics of students' alternative conception. The development procedures, as explained in the methodology, has resulted in 4TMC instrument; however, instrument development is seen as an essential continuous process.

Recommendations

The Based on the results of the study, there are several recommendations for researchers and teachers. For researchers, the findings of this research can be followed up to examine more in how students build their understanding gradually in explaining the concept of particles in substance form changes. The study can be conducted by developing tests that aim to evaluate and diagnose the process of student knowledge formation and development while being able to identify at the level of education where the confusion of understanding occurs. The evaluation becomes more objective, not only reviewed from the student's point of ability but can be reviewed from the teacher's ability. The model of *PCM* multi-faced item response pattern approach becomes one of the important parts recommended for such objectives. In this way, students' ability to develop epistemological knowledge, and their ability to significantly actualize the knowledge gained can be measured well.

On the other hand, for teachers, the results of this study along with the stages of analysis approach used can be a reference in evaluating the progress of learners' learning, as well as determining alternative thinking frameworks of students in explaining the concept of substance change. The information serves as strategic feedback in formulating instructional strategies and preparing remedial learning, especially for students who have difficulty in developing epistemological explanations of substance changes.

Limitations

The limitations of the research are primarily related to the misrepresentation of student reasoning, which may arise in its efforts to connect phenomena and concepts measured in each item. In this context, the student may not excel to explain, because of his incapableness in using his heuristic reasoning. This instrument is not equipped with items that evaluate the heuristic abilities of the student in question. However, researchers decided to record this incompetence as a misconception or vague knowledge. For further research, it is recommended that the instrument be equipped with items that measure students' emotional and heuristic reasoning according to the conceptual framework to be evaluated.

Acknowledgments

The researchers would like to express their gratitude towards the Directorate of Research and Community Service, Ministry of Research and Technology of Republic of Indonesia, for the financial support through the University Basic Research Excellence Grant Program in the Research and Community Service Office of Universitas Negeri Gorontalo, 2020.

References

- Aktan, D. C. (2013). Investigation of students' intermediate conceptual understanding levels: The case of direct current electricity concepts. *European Journal of Physics*, *34*(1), 33–43. https://doi.org/10.1088/0143-0807/34/1/33
- Arslan, H. O., Cigdemoglu, C., & Moseley, C. (2012). A three-tier diagnostic test to assess pre-service teachers' misconceptions about global warming, greenhouse effect, ozone layer depletion, and acid rain. *International Journal of Science Education*, 34(11), 1667–1686. https://doi.org/10.1080/09500693.2012.680618
- Bond, T. G., & Fox, C. M. (2007). *Applying the Rasch Model: Fundamental measurement in the human sciences* (2nd ed.). Routledge.
- Caleon, I. S., & Subramaniam, R. (2010). Do students know what they know and what they don't know? Using a four-tier diagnostic test to assess the nature of students' alternative conceptions. *Research in Science Education*, 40(3), 313–337. https://doi.org/10.1007/s11165-009-9122-4
- Chandrasegaran, A. L., Treagust, D. F., & Mocerino, M. (2007). The development of two tier multiple-choice diagnostic instrument for evaluating secondary school students' ability to describe and explain chemical reactions using multiple levels of representation. *Chemistry Education Research and Practice*, *8*(3), 293–307
- Chi, S., Wang, Z., Luo, M., Yang, Y., & Huang, M. (2018). Student progression on chemical symbol representation abilities at different grade levels (Grades 10–12) across gender. *Chemistry Education Research and Practice*, *19*(4), 1055–1064. https://doi.org/10.1039/c8rp00010g
- Claesgens, J., Scalise, K., Wilson, M., & Stacy, A. (2009). Mapping student understanding in chemistry: The perspectives of chemists. *Science Education*, *93*(1), 56–85. https://doi.org/10.1002/sce.20292.
- Djiks, M. A., Brummer, L., & Kostons, D. (2018). The anonymous reviewer: the relationship between perceived expertise and the perceptions of peer feedback in higher education. *Assessment & Evaluation in Higher Education*, 43(8), 1258-1271. https://doi.org/10.1080/02602938.2018.1447645
- Duncan, R. G., & Hmelo-Silver, C. E. (2009). Learning progressions: Aligning curriculum, instruction, and assessment. *Journal of Research in Science Teaching*, 46(6), 606–609. https://doi.org/10.1002/tea.20316
- Duschl, R., Maeng, S., & Sezen, A. (2011). Learning progressions and teaching sequences: A review and analysis. *Studies in Science Education*, 47(2), 123–182. https://doi.org/10.1080/03057267.2011.604476
- Emden, M., Weber, K., & Sumfleth, E. (2018). Evaluating a learning progression on "Transformation of Matter" on the lower secondary level. *Chemistry Education Research and Practice*, *19*(4), 1096–1116. https://doi.org/10.1039/c8rp00137e
- Goulas, S., & Megalokonomou, R. (2021). Knowing who you actually are: The effect of feedback on short-and longerterm outcomes. *Journal of Economic Behavior & Organization, 183*, 589-615. https://doi.org/10.1016/j.jebo.2021.01.013
- Guo, M., & Leung, F. K. S. (2021). Achievement goal orientations, learning strategies, and mathematics achievement: A comparison of Chinese Miao and Han students. *Psychology in the Schools.* 58(1), 107-123.

https://doi.org10.1002/pits.22424

- Haarala-Muhonen, A., Ruohoniemi, M., Parpala, A., Komulainen, E., & Lindblom-Ylänne, S. (2016). How do the different study profiles of first-year students predict their study success, study progress and the completion of degrees? *Higher Education*, 74(6), 949–962. https://doi.org/ 10.1007/s10734-016-0087-8
- Habiddin, & Page, E. M. (2019). Development and validation of a four-tier diagnostic instrument for chemical kinetics (FTDICK). *Indonesian Journal of Chemistry*, *19*(3), 720–736. https://doi.org/10.22146/ijc.39218
- Hadenfeldt, J. C., Bernholt, S., Liu, X., Neumann, K., & Parchmann, I. (2013). Using ordered multiple-choice items to assess students' understanding of the structure and composition of matter. *Journal of Chemical Education*, *90*(12), 1602–1608. https://doi.org/10.1021/ed3006192
- Hasan, S., Bagayoko, D., & Kelley, E. L. (1999). Misconceptions and the certainty of response index (CRI). *Physics Education*, 34(5), 294–299. https://doi.org/10.1088/0031-9120/34/5/304
- Herrmann-Abell, C. F., & Deboer, G. E. (2016). Using rasch modeling and option probability curves to diagnose students' misconceptions. *American Educational Research Association*, 8(12), 1–12
- Hoe, K. Y., & Subramaniam, R. (2016). On the prevalence of alternative conceptions on acid-base chemistry among secondary students: Insights from cognitive and confidence measures. *Chemistry Education Research and Practice*, 17(2), 263–282. https://doi.org/10.1039/c5rp00146c
- Jin, H., Mikeska, J. N., Hokayem, H., & Mavronikolas, E. (2019). Toward coherence in curriculum, instruction, and assessment: A review of learning progression literature. *Science Education*, 103(5), 1206–1234. https://doi.org/10.1002/sce.21525
- Karagiannopoulou, E., Milienos, F. S., & Rentzios, C. (2020). Grouping learning approaches and emotional factors to predict students' academic progress. *International Journal of School & Educational Psychology*, 9(1), 1– 18. https://doi.org.10.1080/2168363.2020.183241
- Klassen, S. (2006). Contextual assessment in science education: Background, issues, and policy. *Science Education*, *90*(5), 820–851. https://doi.org/10.1002/sce.20150
- Latifi, S., Noroozi, O., & Talaee, E. (2021). Peer feedback or peer feedforward? Enhancing students' argumentative peer learning processes and outcomes. *British Journal of Educational Technology*, 52(2), 768-784. https://doi.org/10.1111/bjet.13054
- Lee, K., & Keller, J. M. (2021). Use of the ARCS model in education: A literature review. *Computers & Education*, 122(1), 54-62. https://doi.org/10.1016/j.compedu.2018.03.019
- Lin, P. Y., Chai, C. S., Jong, M. S. Y., Dai, Y., Guo, Y., & Qin, J. (2021). Modeling the structural relationship among primary students' motivation to learn artificial intelligence. *Computers and Education: Artificial Intelligence*, *2*(1), 1-7
- Laliyo, Botutihe, & Panigoro. (2019). The development of two-tier instrument based on distractor to assess conceptual understanding level and student misconceptions in explaining redox reactions. *International Journal of Learning, Teaching and Educational Research*, *18*(9), 216–237. https://doi.org/10.26803/ijlter.18.9.12
- Linacre, J. M. (2012). A user's guide to WINSTEPS® MINISTEP Rasch-model computer program: Program manual 3.75.0. winsteps.com.
- Linacre, J. M. (2020). A User's Guide to WINSTEPS ® MINISTEP Rasch-Model Computer Programs Program Manual 4.5.1. winsteps.com.
- Ling Lee, W., Chinna, K., & Sumintono, B. (2020). Psychometrics assessment of HeartQoL questionnaire: A Rasch analysis. *European Journal of Preventive Cardiology*. Advance online publication. https://doi.org/10.1177/2047487320902322
- Liu, X. (2012). Developing measurement instruments for science education research. In B. Fraser, K. G. Tobin, & C. J. McRobbie (Eds.), *Second international handbook of science education* (pp. 651–665). Springer Netherlands
- Lu, S., & Bi, H. (2016). Development of a measurement instrument to assess students' electrolyte conceptual understanding. *Chemistry Education Research and Practice*, *17*(4), 1030–1040. https://doi.org/10.1039/c6rp00137h
- Morell, L., Collier, T., Black, P., & Wilson, M. (2017). A construct-modeling approach to develop a learning progression of how students understand the structure of matter. *Journal of Research in Science Teaching*, *54*(8), 1024–1048. https://doi.org/10.1002/tea.21397
- Neumann, K., Viering, T., Boone, W. J., & Fischer, H. E. (2013). Towards a learning progression of energy. *Journal of Research in Science Teaching*, 50(2), 162–188. https://doi.org/10.1002/tea.21061

840 | LALIYO ET AL. / Implementation of Four-Tier Instruments Based on the Partial Credit Model

- Park, M., Liu, X., & Waight, N. (2017). Development of the connected chemistry as formative assessment pedagogy for high school chemistry teaching. *Journal of Chemical Education*, 94(3), 273–281. https://doi.org/10.1021/acs.jchemed.6b00299
- Peterson, R. F., Treagust, D. F., & Garnett, P. (1989). Development and application of a diagnostic instrument to evaluate grade-11 and -12 students' concepts of covalent bonding and structure following a course of instruction. *Journal of Research in Science Teaching*, *26*(4), 301–314. https://doi.org/10.1002/tea.3660260404
- Rogat, A. (2011). *Developing learning progressions in support of the new science standards: A RAPID workshop series.* CPRE Research Reports. http://repository.upenn.edu/cpre_researchreports/66
- Smith, C. L., Wiser, M., Anderson, C. W., & Krajcik, J. (2006). Implications of research on children's learning for standards and assessment: A proposed learning progression for matter and the atomic-molecular theory. *Measurement: Interdisciplinary Research & Perspective*, *4*(1–2), 1–98. https://doi.org/10.1080/15366367.2006.9678570
- Sumintono, B., & Widhiarso, W. (2014). *Aplikasi model Rasch untuk penelitian ilmu-ilmu sosial* [Application of Rasch model in social science research]. Trim Komunikata.
- Sutiani, A., Situmorang, M., & Silalahi, A. (2021). Implementation of an Inquiry Learning Model with Science Literacy to Improve Student Critical Thinking Skills. *International Journal of Instruction*, *14*(2), 117-138.
- Sreenivasulu, B., & Subramaniam, R. (2013). University students' understanding of chemical thermodynamics. *International Journal of Science Education*, *35*(4), 601-635.
- Testa, I., Capasso, G., Colantonio, A., Galano, S., Marzoli, I., Scotti di Uccio, U., Trani, F., & Zappia, A. (2019). Development and validation of a university students' progression in learning quantum mechanics through exploratory factor analysis and Rasch analysis. *International Journal of Science Education*, 41(3), 388–417. https://doi.org/10.1080/09500693.2018.1556414
- Treagust, D. F. (1988). Development and use of diagnostic tests to evaluate students' misconceptions in science. *International Journal of Science Education*, *10*(2), 159–169. https://doi.org/10.1080/0950069880100204
- Tyson, L., Treagust, D. F., & Bucat, R. B. (1999). The complexity of teaching and learning chemical equilibrium. *Journal of Chemical Education*, 76(2-4), 554–558. https://doi.org/10.1021/ed077p1560.1
- Wilson, M. (2005). *Constructing measures: An item response modeling approach*. Lawrence Erlbaum Associates, Inc. https://doi.org/10.4324/9781410611697
- Wilson, M. (2008). Cognitive diagnosis using item response models. *Journal of Psychology/ Zeitschrift Für Psychologie*, 216(2), 74–88. https://doi.org/10.1027/0044-3409.216.2.74
- Wilson, M. (2009). Measuring progressions: Assessment structures underlying a learning progression. *Journal of Research in Science Teaching*, 46(6), 716–730. https://doi.org/10.1002/tea.20318
- Wilson, M. (2012). Responding to a challenge that learning progressions pose to measurement practice. In A. C. Alonzo & A. W. Gotwals (Eds.), *Learning progression in science* (pp. 317–344). Sense Publishers. https://doi.org/10.1007/978-94-6091-824-7