



International Conference on Artificial Intelligence in Education Technology (AIET 2021)

TIM SCHLIPPE & JÖRG SAWATZKI

**CROSS-LINGUAL**

**AUTOMATIC SHORT ANSWER GRADING**

Wuhan, China

July 4, 2021

# AGENDA

---

**Motivation**

---

**1**

---

**Cross-Lingual Automatic Short Answer Grading**

---

**2**

---

**Experimental Setup**

---

**3**

---

**Experiments and Results**

---

**4**

---

**Conclusion and Future Work**

---

**5**

# MOTIVATION

# MOTIVATION

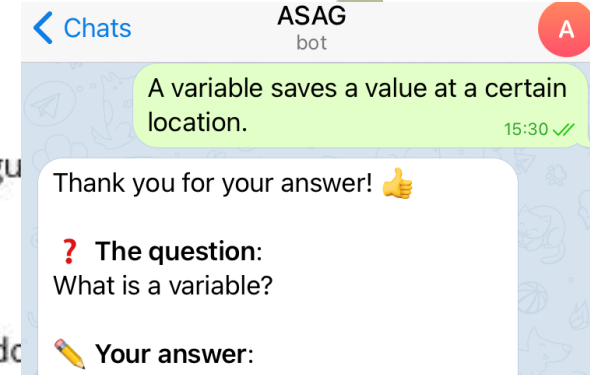


# MOTIVATION

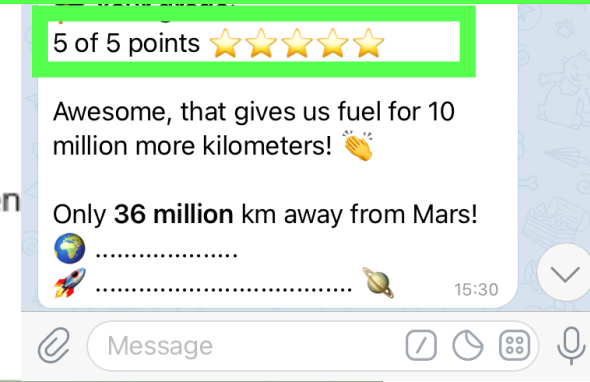


Sources: United Nations: Sustainable Development Goals: 17 Goals to Transform our World (2021); OpenClipart-Vectors/154119/Pixabay; Statista: The Most Spoken Languages Worldwide in 2019 (2020).

# MOTIVATION



## CROSS-LINGUAL AUTOMATIC SHORT ANSWER GRADING



# MOTIVATION



<b>Question</b>	What is a variable?
<b>Model answer</b>	A location in memory that can store a value.
<b>Answer 1</b>	Eine Stelle im Speicher, die einen Wert speichern kann.
<b>Grading: Answer 1</b>	<b>5 of 5 points</b>
<b>Answer 2</b>	变量可以是整数，也可以是程序中的字符串。
<b>Grading: Answer 2</b>	<b>2 of 5 points</b>

# 02

## **CROSS-LINGUAL AUTOMATIC SHORT ANSWER GRADING**



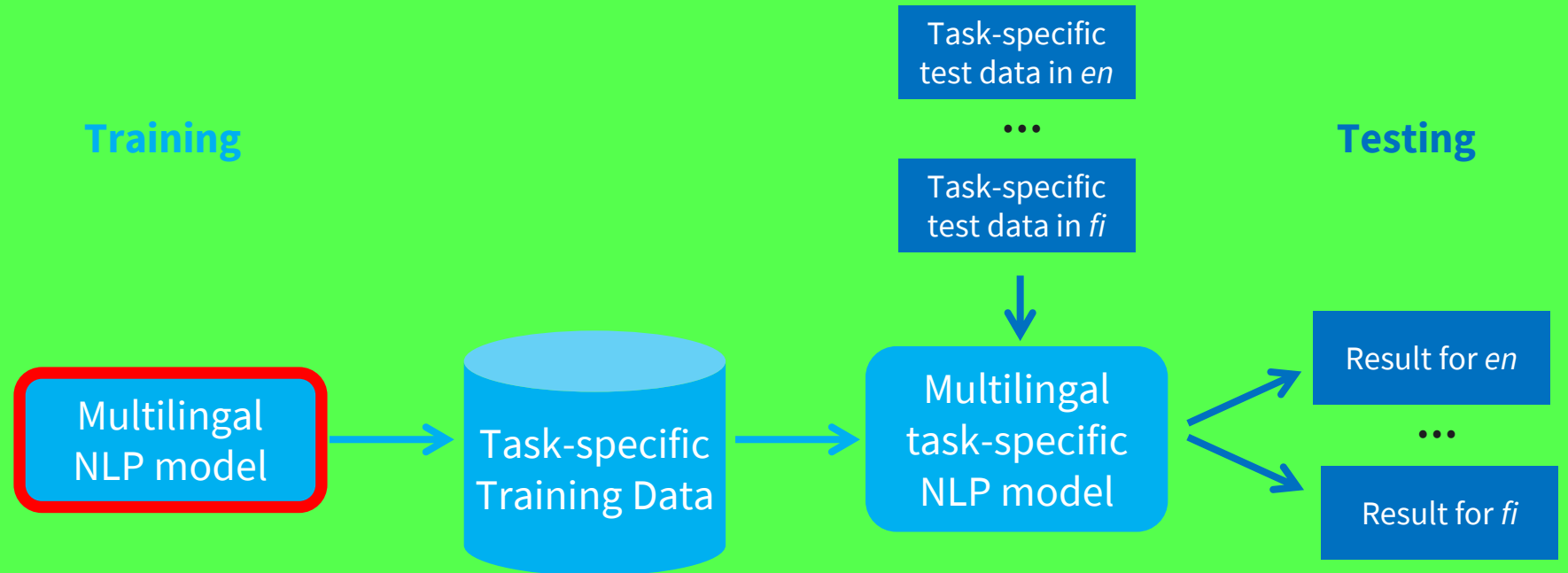
# CROSS-LINGUAL AUTOMATIC SHORT ANSWER GRADING

## MULTILINGUAL NLP MODELS

**Transfer learning**  
**Cross-lingual transfer**

e.g., Multilingual BERT  
(Devlin et al., 2019;  
Pires et al., 2019),  
RoBERTa (Liu et al., 2019),  
XLM-R (Conneau, 2018)

**Training**



# CROSS-LINGUAL AUTOMATIC SHORT ANSWER GRADING

## MULTILINGUAL NLP MODELS

Transfer learning  
Cross-lingual transfer

e.g., Multilingual BERT

### Multilingual BERT

PLIES et al., 2019),  
RoBERTa (Liu et al., 2019),  
XLM-R (Conneau, 2018)

Training

Multilingual  
NLP model

Task-specific  
Training Data

Task-specific  
test data in *en*

...

Task-specific  
test data in *fi*

Testing

Multilingual  
task-specific  
NLP model

Result for *en*

...

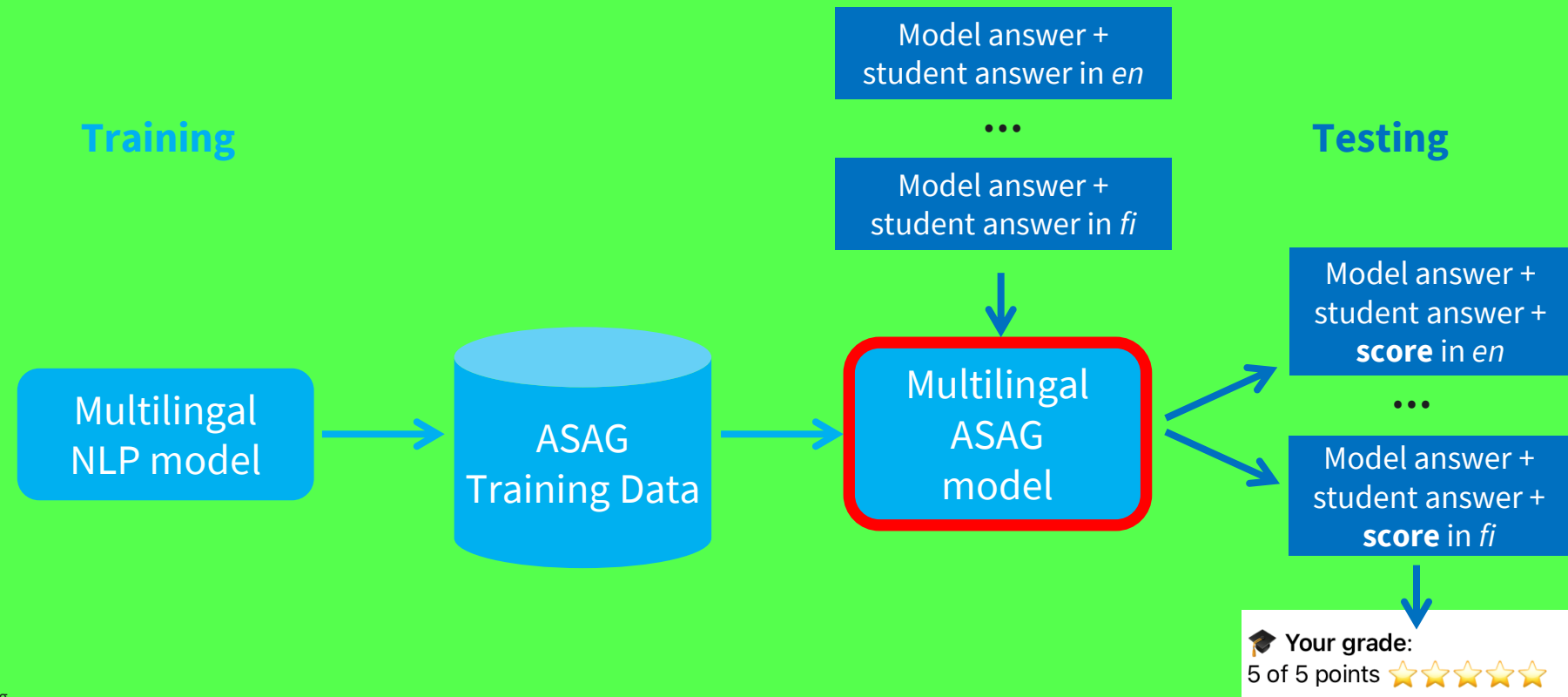
Result for *fi*

# CROSS-LINGUAL AUTOMATIC SHORT ANSWER GRADING

## AUTOMATIC SHORT ANSWER GRADING

### Deep learning

e.g., (Burrows et al., 2014;  
Camus & Filighera, 2020;  
Sawatzki et al., 2021;  
Schlippe & Sawatzki, 2021b)

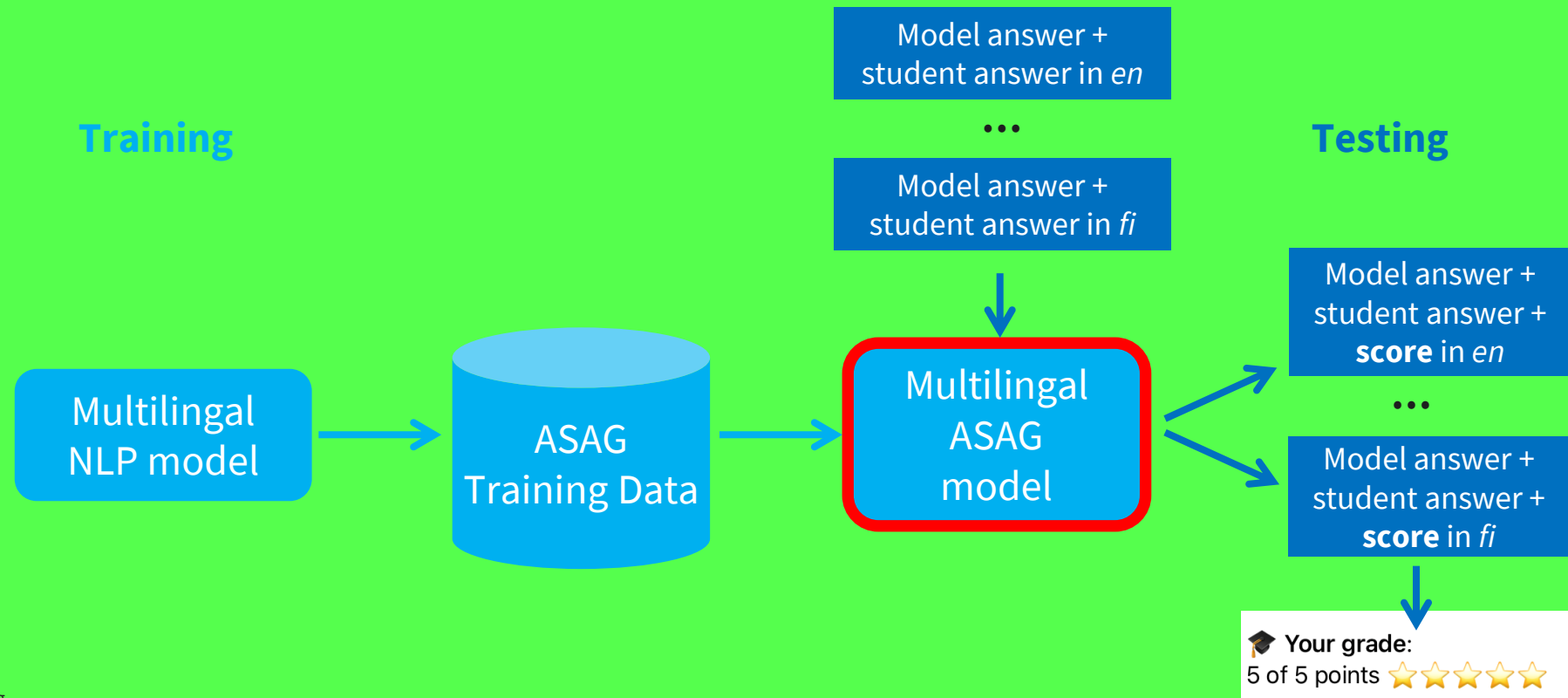


# CROSS-LINGUAL AUTOMATIC SHORT ANSWER GRADING

## AUTOMATIC SHORT ANSWER GRADING

Deep learning

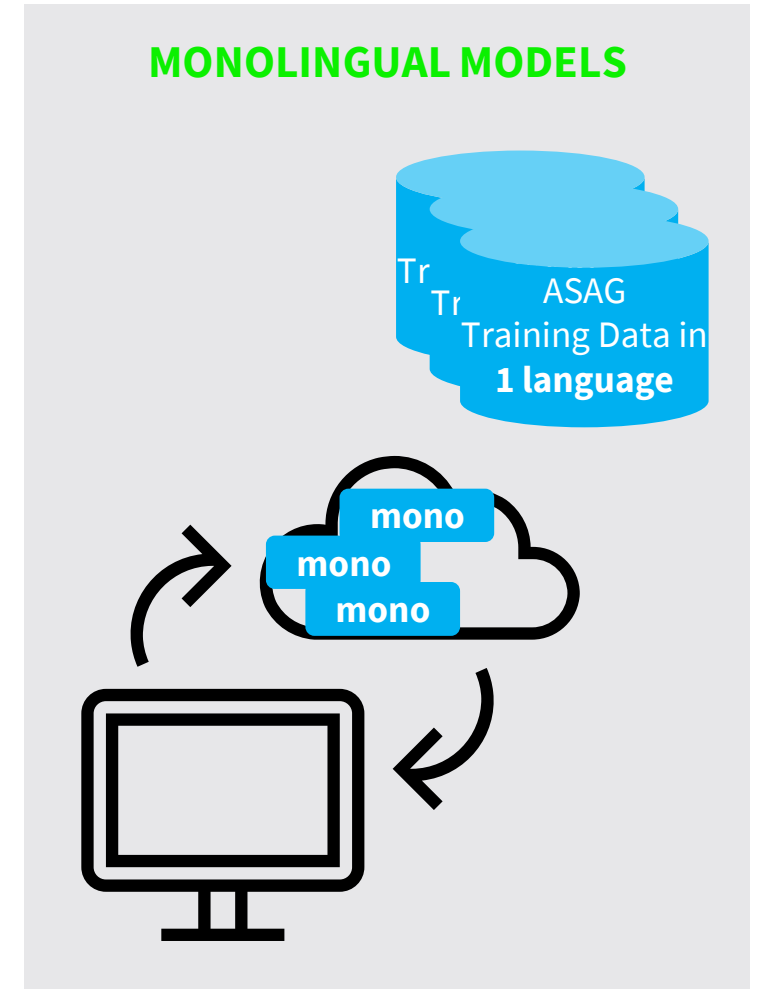
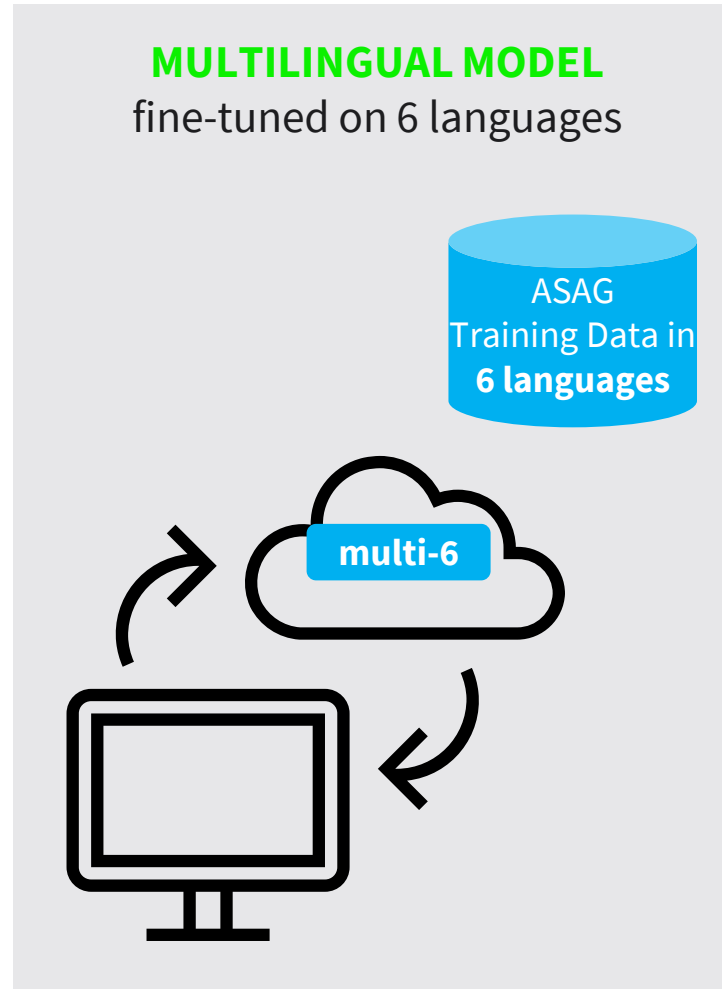
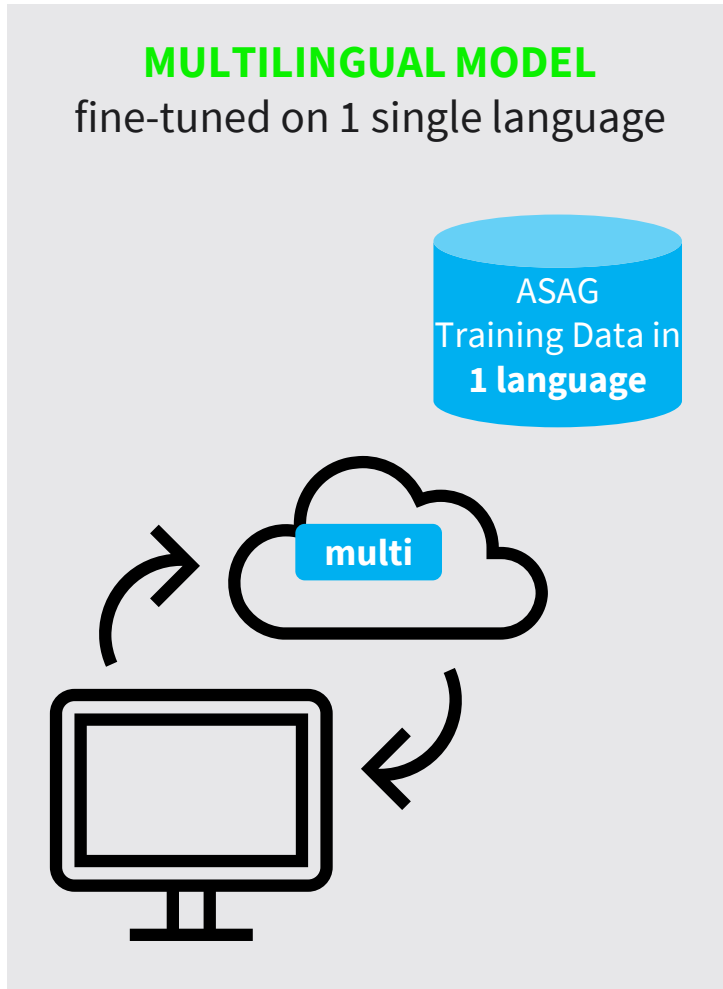
## Cross-lingual Automatic Short Answer Grading



# 03

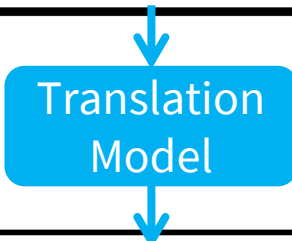
## EXPERIMENTAL SETUP

# EXPERIMENTAL SETUP

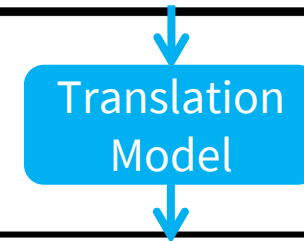


# EXPERIMENTAL SETUP

<b>Question</b>	What is a variable?
<b>Model answer</b>	A location in memory that can store a value.
<b>Example: Answer 1</b>	A variable is a location in memory where a value can be stored.
<b>Grading: Answer 1</b>	5 of 5 points
<b>Example: Answer 2</b>	Variable can be an integer or a string in a program.
<b>Grading: Answer 2</b>	2 of 5 points



**Google's Neural Machine Translation System**  
(Wu et al., 2016; Aiken, 2019)



<b>Question</b>	Was ist eine Variable?
<b>Model answer</b>	Eine Stelle im Speicher, die einen Wert speichern kann.
<b>Example: Answer 1</b>	Eine Variable ist ein Ort im Speicher, an dem ein Wert gespeichert werden kann.
<b>Grading: Answer 1</b>	5 of 5 points
<b>Example: Answer 2</b>	Eine Variable kann in einem Programm ein Integer oder ein String sein.
<b>Grading: Answer 2</b>	2 of 5 points

<b>Question</b>	什么是变量？
<b>Model answer</b>	内存中可以存储值的位置。
<b>Example: Answer 1</b>	变量是内存中可以存储值的位置。
<b>Grading: Answer 1</b>	5 of 5 points
<b>Example: Answer 2</b>	变量可以是整数，也可以是程序中的字符串。
<b>Grading: Answer 2</b>	2 of 5 points

# 04

## EXPERIMENTS AND RESULTS



# EXPERIMENTS AND RESULTS

	multi+ en	multi+ de	multi+ nl	multi+ jp	multi+ zh	multi+ fi	multi+ 6	mono
en	0.45	0.61	0.64	0.68	0.63	0.63	<b>0.43</b>	0.43
ceb	0.70	0.73	0.72	0.68	0.72	0.71	<b>0.63</b>	-
sv	0.63	0.67	0.68	0.73	0.72	0.68	<b>0.48</b>	-
de	0.64	0.51	0.67	0.70	0.70	0.65	0.46	<b>0.45</b>
fr	0.61	0.66	0.64	0.67	0.70	0.67	<b>0.54</b>	-
nl	0.62	0.64	0.52	0.70	0.73	0.67	<b>0.45</b>	0.47
ru	0.68	0.73	0.83	0.74	0.75	0.78	<b>0.52</b>	-
it	0.62	0.65	0.72	0.71	0.73	0.70	<b>0.52</b>	-
es	0.61	0.68	0.76	0.68	0.72	0.65	<b>0.49</b>	-
pl	0.62	0.71	0.77	0.69	0.72	0.68	<b>0.51</b>	-
vi	0.71	0.72	0.84	0.77	0.73	0.71	<b>0.52</b>	-
jp	0.66	0.70	0.73	0.49	0.63	0.71	<b>0.44</b>	0.53
zh	0.63	0.71	0.77	0.69	0.50	0.79	<b>0.41</b>	0.44
ar	0.72	0.78	0.85	0.78	0.76	0.76	<b>0.59</b>	-
uk	0.65	0.70	0.82	0.73	0.73	0.75	<b>0.54</b>	-
pt	0.59	0.67	0.75	0.69	0.73	0.69	<b>0.50</b>	-
fa	0.64	0.66	0.71	0.67	0.70	0.69	<b>0.56</b>	-
ca	0.64	0.70	0.74	0.70	0.76	0.67	<b>0.53</b>	-
sr	0.69	0.81	0.83	0.76	0.79	0.86	<b>0.56</b>	-
id	0.66	0.68	0.69	0.70	0.79	0.63	<b>0.49</b>	-
no	0.63	0.69	0.65	0.75	0.71	0.69	<b>0.45</b>	-
ko	0.70	0.70	0.76	0.66	0.66	0.67	<b>0.58</b>	-
fi	0.69	0.79	0.77	0.77	0.73	0.52	0.47	<b>0.45</b>
hu	0.69	0.76	0.81	0.72	0.76	0.69	<b>0.54</b>	-
cs	0.62	0.77	0.82	0.72	0.78	0.71	<b>0.51</b>	-
sh	0.66	0.77	0.79	0.74	0.78	0.79	<b>0.53</b>	-

# EXPERIMENTS AND RESULTS

**0.75 POINTS  
HUMAN GRADER VARIABILITY**

	multi+ en	multi+ de	multi+ nl	multi+ jp	multi+ zh	multi+ fi	multi+ 6	mono
en	0.45	0.61	0.64	0.68	0.63	0.63	<b>0.43</b>	0.43
ceb	0.70	0.73	0.72	0.68	0.72	0.71	<b>0.63</b>	-
sv	0.63	0.67	0.68	0.73	0.72	0.68	<b>0.48</b>	-
de	0.64	0.51	0.67	0.70	0.70	0.65	0.46	<b>0.45</b>
fr	0.61	0.66	0.64	0.67	0.70	0.67	<b>0.54</b>	-
nl	0.62	0.64	0.52	0.70	0.73	0.67	<b>0.45</b>	0.47
ru	0.68	0.73	0.83	0.74	0.75	0.78	<b>0.52</b>	-
it	0.62	0.65	0.72	0.71	0.73	0.70	<b>0.52</b>	-
es	0.61	0.68	0.76	0.68	0.72	0.65	<b>0.49</b>	-
pl	0.62	0.71	0.77	0.69	0.72	0.68	<b>0.51</b>	-
vi	0.71	0.72	0.84	0.77	0.73	0.71	<b>0.52</b>	-
jp	0.66	0.70	0.73	0.49	0.63	0.71	<b>0.44</b>	0.53
zh	0.63	0.71	0.77	0.69	0.50	0.79	<b>0.41</b>	0.44
ar	0.72	0.78	0.85	0.78	0.76	0.76	<b>0.59</b>	-
uk	0.65	0.70	0.82	0.73	0.73	0.75	<b>0.54</b>	-
pt	0.59	0.67	0.75	0.69	0.73	0.69	<b>0.50</b>	-
fa	0.64	0.66	0.71	0.67	0.70	0.69	<b>0.56</b>	-
ca	0.64	0.70	0.74	0.70	0.76	0.67	<b>0.53</b>	-
sr	0.69	0.81	0.83	0.76	0.79	0.86	<b>0.56</b>	-
id	0.66	0.68	0.69	0.70	0.79	0.63	<b>0.49</b>	-
no	0.63	0.69	0.65	0.75	0.71	0.69	<b>0.45</b>	-
ko	0.70	0.70	0.76	0.66	0.66	0.67	<b>0.58</b>	-
fi	0.69	0.79	0.77	0.77	0.73	0.52	0.47	<b>0.45</b>
hu	0.69	0.76	0.81	0.72	0.76	0.69	<b>0.54</b>	-
cs	0.62	0.77	0.82	0.72	0.78	0.71	<b>0.51</b>	-
sh	0.66	0.77	0.79	0.74	0.78	0.79	<b>0.53</b>	-

Mean Absolute Error  
out of 5 points

# EXPERIMENTS AND RESULTS

	<b>multi+ en</b>	<b>multi+ 6</b>	<b>rel. improvement</b>
<b>en</b>	0.45	0.43	4.4%
<b>de</b>	0.64	0.46	28.1%
<b>nl</b>	0.62	0.45	27.4%
<b>jp</b>	0.66	0.44	33.3%
<b>zh</b>	0.63	0.41	<b>34.9%</b>
<b>fi</b>	0.69	0.47	31.9%

Mean Absolute Error  
out of 5 points  
(Human grader variability: 0.75)

→ Up to 35% improvement by adding more languages

# EXPERIMENTS AND RESULTS

	<b>multi+</b> <i>L<sub>target</sub></i>	<b>multi+</b> <b>6</b>	<b>rel.</b> <b>improvement</b>
<b>en</b>	0.45	0.43	4.4%
<b>de</b>	0.51	0.46	9.8%
<b>nl</b>	0.52	0.45	13.5%
<b>jp</b>	0.49	0.44	10.2%
<b>zh</b>	0.50	0.41	<b>18.0%</b>
<b>fi</b>	0.52	0.47	9.6%

Mean Absolute Error  
out of 5 points

(Human grader variability: 0.75)

→ Even with data in the target language, adding the 5 languages provides improvements of up to 18%

# 05

## CONCLUSION AND FUTURE WORK

# CONCLUSION AND FUTURE WORK

## Conclusion

- Potential of cross-lingual automatic short answer grading
- Analysis on 26 languages
- Mean Absolute Errors (MER) between 0.41 and 0.72 points out of 5 points
- Less discrepancy than 2 human graders
- Augmenting training data with machine translated task-specific data for fine-tuning improves multilingual models
- Results experimentally

# CONCLUSION AND FUTURE WORK

## Conclusion

- Potential of cross-lingual automatic short answer grading
- Analysis on 26 languages
- Mean Absolute Errors (MER) between 0.41 and 0.72 points out of 5 points
- Less discrepancy than 2 human graders
- Augmenting training data with machine translated task-specific data for fine-tuning improves multilingual models
- Results experimentally

## Future Work

- Extension to other languages
- Integration and application for online exams
- Interactive training programs for exam preparation (Schlippe & Sawatzki, 2021a)
- Address the issue of explainability to provide better support to human graders

# FUTURE WORK

Question 1/3: What are ruminants? **Question to be answered**

**Student answer** **Given answer**  
An example are cows that regurgitate their nourishment and chew it up again. They do this in order to be able to digest food better.

**Model answer (Maximum 2 points)** **Model answer with score**  
Ruminants regurgitate predigested food from their stomachs (1 point) and chew it again (1 point).

AI Assistant

AI PREDICTION	EXPLANATION
<b>2 points</b>	Highlighted in color, you can see the individual words were for the
An example are cows that regurgitate their nourishment it up again. They do this in order to be able to digest f	Information provided by the AI assistant. Different information is available here for each concept.
Irrelevant <input type="checkbox"/> <input checked="" type="checkbox"/> Very relevant	

Your scoring (in points)

2

**Your task is to assign a score (in points) using the model answer and the AI assistants.**

<https://tinyurl.com/AI-grade>

**Try it out and  
support us.**



**THANK YOU**

Tim Schlippe

 [tim.schlippe@iu.org](mailto:tim.schlippe@iu.org)

# REFERENCES

## Literature

- **United Nations: Sustainable Development Goals: 17 Goals to Transform our World (2021):**  
<https://www.un.org/sustainabledevelopment/sustainable-development-goals>
- **Statista: The Most Spoken Languages Worldwide in 2019 (2020):**  
<https://www.statista.com/statistics/266808/the-most-spoken-languages-worldwide>
- **Schlippe, T. & Sawatzki, J. (2021a):** *AI-based Multilingual Interactive Exam Preparation*. In: Proceedings of The Learning Ideas Conference 2021 (14th annual conference), ALICE - Special Conference Track on Adaptive Learning via Interactive, Collaborative and Emotional Approaches, New York, New York, 14-18 June 2021. In: *Advances in Intelligent Systems and Computing*, Springer, 2021.
- **Devlin, J., Chang, M.W., Lee, K., & Toutanova, K. (2019):** *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. In: NAACL-HLT, Minneapolis, Minnesota.
- **Pires, T., Schlinger, E., & Garrette, D. (2019):** *How Multilingual is Multilingual BERT?* In: ACL. Florence, Italy.
- **Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019):** *RoBERTa: A Robustly Optimized BERT Pretraining Approach*. arXiv:1907.11692.

## Images

- **Images provided by OpenClipart-Vectors/154119/Pixabay.**  
(<https://pixabay.com/vectors/international-project-world-154119> [last access: 16.16.2021])
- **Coursebook (NLP and Computer Vision):** IU International University of Applied Sciences. Version No.: 001-2020-1211.

## Literature

- **Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., & Stoyanov, V. (2018):** *Unsupervised Cross-lingual Representation Learning at Scale*. arXiv:1911.02116.
- **Burrows, S., Gurevych, I., & Stein, B. (2014):** *The Eras and Trends of Automatic Short Answer Grading*. In: IJAIED 25, 60–117.
- **Camus, L. & Filighera, A. (2020):** *Investigating Transformers for Automatic Short Answer Grading*. AIED. Cyberspace.
- **Sawatzki, J., Schlippe, T., & Benner-Wickner, M. (2021):** *Deep Learning Techniques for Automatic Short Answer Grading: Predicting Scores for English and German Answers*. In: AIET, Wuhan, China.
- **Schlippe, T. & Sawatzki, J. (2021b):** *Cross-Lingual Automatic Short Answer Grading*. In: AIET, Wuhan, China.
- **Mohler, M., Bunescu, R., & Mihalcea, R. (2011):** *Learning to Grade Short Answer Questions using Semantic Similarity Measures and Dependency Graph Alignments*. In: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies. pp. 752–762. Association for Computational Linguistics, Portland, Oregon, USA.