



Hochschule  
Bonn-Rhein-Sieg  
University of Applied Sciences



# Analysis of Active Learning Mechanism Applied to Language Models for Computer Assisted Short Answer Grading

September 26, 2022

Elanton Fernandes

*Advisors*

Prof. Dr. Paul G. Plöger, M.Sc Tim Metzler

# Table of Contents

1. Motivation
2. Problem Statement
3. State of the Art
4. Approach
5. Evaluation
6. Results
7. Summary
8. Future Work
9. Extra Slides



# Motivation

In universities with an increase in number of student every semester, the number of tests conducted also increases. This means that:

- The professor spends more time in correcting student exams than preparing for lectures.
- If students are not assigned full scores for on a test, they expect a meaningful feedback from the professor.

# Motivation

Consider the following dummy scenario:

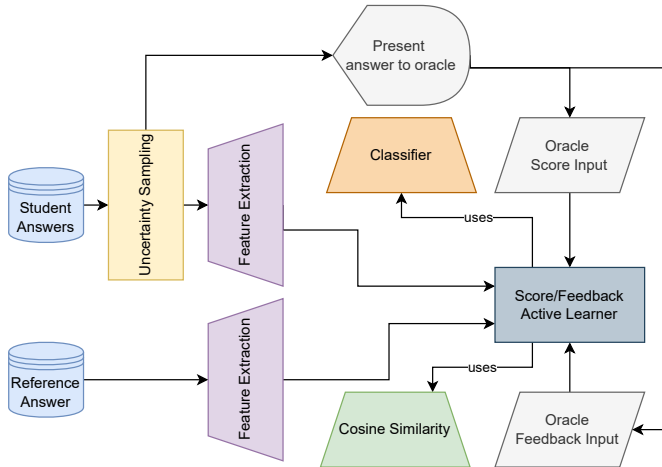
- 80 students enrolled in a class.
- Tests are conducted bi-weekly.
- Professor requires 15 minutes to evaluate one student test.
- Total time spent by the professor to evaluate all tests per week is 10 hours.

# Problem Statement

- To automate the evaluation of student tests while still keeping the oracle/professor in the loop.
- Allow the assignment of meaningful feedback to student answers indicating their mistakes.

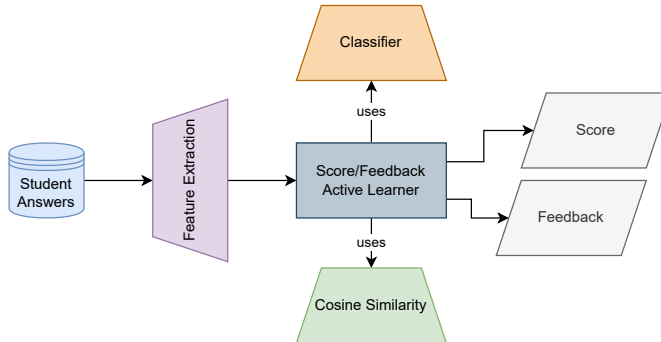
# Approach

## Training cycle



# Approach

## *Prediction cycle*



# Approach

## Uncertainty Sampling

Uncertainty sampling is a query strategy that queries the instances about which it is least certain how to label. We use uncertainty sampling variant might query the instance whose prediction is the least confident:

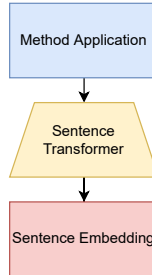
$$x_{LC} = \operatorname{argmin}_x P(\hat{y}|x; \theta) \quad (1)$$

Where  $x$  is the feature,  $y$  is the class label prediction, and  $\hat{y} = \operatorname{argmax}_y P(y|x; \theta)$  is the class label that has the largest posterior probability using model  $\theta$ .



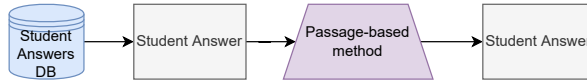
# Approach

## *Feature Extraction: Overview*



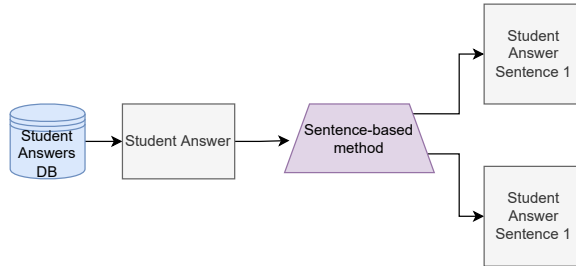
# Approach

## *Feature Extraction: Passage-based method*



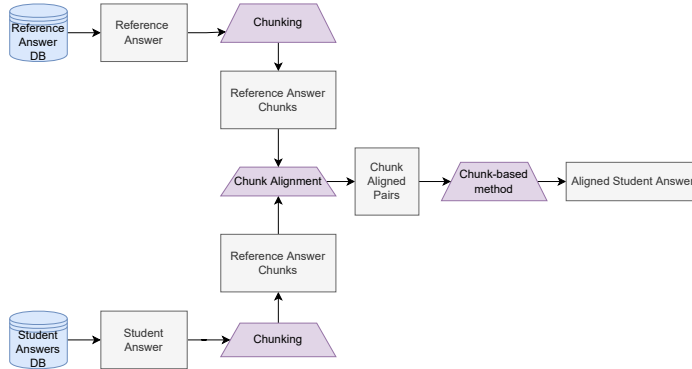
# Approach

## *Feature Extraction: Sentence-based method*



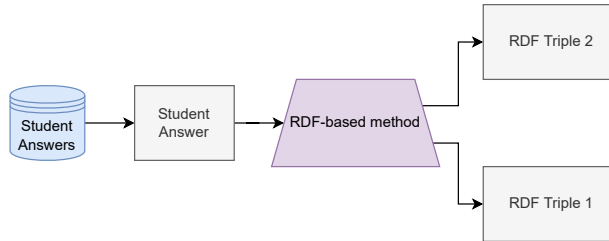
# Approach

## *Feature Extraction: Chunk-based method*



# Approach

## *Feature Extraction: RDF-based method*



# Approach

## Language Models

Model:	Base model	Number Training tuples
all-mpnet-base-v2[?]	microsoft/mpnet-base.	1.17B
all-distilroberta-v1[?]	distilroberta-base	1.12B
all-MiniLM-L12-v2[?]	microsoft/MiniLM-L12-H384-uncased	1.17B
multi-qa-distilbert-cos-v1[?]	distilbert-base	214M
all-MiniLM-L6-v2[?]	nreimers/MiniLM-L6-H384-uncased	1.17B

**Table 1:** Displays pre-trained language models with their base model used in training and number of training tuples used[?].

# Evaluation

## Score

- Pearsons Correlation

$$\rho(y, \hat{y}) = \frac{\text{cov}(y, \hat{y})}{\sigma_y \sigma_{\hat{y}}} \quad (2)$$

- RMSE Score

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (3)$$

Where  $y$  represents actual grade and  $\hat{y}$  represents predicted grade with  $\sigma_y$  and  $\sigma_{\hat{y}}$  computed as the standard deviation of  $y$  and  $\hat{y}$

# Evaluation

## Feedback

Question	What is a variable?
Reference Answer	A location in memory that can store a value.
Student Answer	a value/word that can assume any of a set of values
Feedback A	correct
Feedback B	missing keywords: location in memory
Feedback C	A variable is a location in memory that stores a value

Table 2: Presented survey to participants.

Agreement Score = Participants agreed with most rated feedback/100



# Results

Score: Pearson Correlation (Methods)

Dataset	M1	M2	M3	M4
Mohler []	<b>0.826</b>	<b>0.791</b>	<b>0.816</b>	0.782
NN Exam []	<b>0.941</b>	<b>0.828</b>	0.561	<b>0.846</b>
AMR Exam []	<b>0.658</b>	0.458	<b>0.640</b>	<b>0.428</b>

(a)

Dataset	M1	M2	M3	M4
Mohler []	0.689	0.627	0.687	<b>0.792</b>
NN Exam []	0.889	0.791	<b>0.638</b>	0.664
AMR Exam []	0.622	<b>0.474</b>	0.593	<b>0.428</b>

(b)

**Table 3:** Comparison of Pearson Correlation between Random Forest (a) and AdaBoost (b) classifiers. Where M1: Passage-based, M2: Sentence-based, M3:Chunk-based, and M4: RDF-based method.

# Results

Score: Pearson Correlation (Language Models)

Dataset	LM1	LM2	LM3	LM4	LM5
Mohler []	<b>0.802</b>	<b>0.797</b>	<b>0.796</b>	<b>0.796</b>	<b>0.789</b>
NN Exam []	<b>0.732</b>	<b>0.670</b>	<b>0.705</b>	<b>0.755</b>	<b>0.760</b>
AMR Exam []	0.453	<b>0.518</b>	<b>0.525</b>	<b>0.523</b>	<b>0.503</b>

(a)

Dataset	LM1	LM2	LM3	LM4	LM5
Mohler []	0.659	0.673	0.211	0.544	0.499
NN Exam []	0.614	0.653	0.704	0.698	0.605
AMR Exam []	<b>0.502</b>	0.440	0.430	0.508	0.467

(b)

Table 4: Comparison of Pearson Correlation between Random Forest (a) and AdaBoost (b) classifiers with language models (LM).

# Results

Score: Root Mean Square Error (Methods)

Dataset	M1	M2	M3	M4
Mohler []	<b>0.893</b>	<b>0.949</b>	<b>0.920</b>	0.942
NN Exam []	<b>0.296</b>	<b>0.520</b>	<b>0.433</b>	<b>0.522</b>
AMR Exam []	<b>0.596</b>	0.716	<b>0.596</b>	<b>0.736</b>

(a)

Dataset	M1	M2	M3	M4
Mohler []	1.218	1.226	1.169	<b>0.920</b>
NN Exam []	0.405	0.571	0.495	0.741
AMR Exam []	0.616	<b>0.707</b>	0.630	0.741

(b)

Table 5: Comparison of RMSE score between Random Forest (a) and AdaBoost (b) classifiers with methods (M).

# Results

*Score: Root Mean Square Error (Language Models)*

Dataset	LM1	LM2	LM3	LM4	LM5
Mohler []	<b>0.931</b>	<b>0.941</b>	<b>0.941</b>	<b>0.941</b>	<b>0.956</b>
NN Exam []	<b>0.484</b>	0.591	<b>0.558</b>	<b>0.490</b>	<b>0.492</b>
AMR Exam []	0.735	<b>0.680</b>	<b>0.676</b>	0.684	<b>0.698</b>

(a)

Dataset	LM1	LM2	LM3	LM4	LM5
Mohler []	1.182	1.163	1.667	1.278	1.363
NN Exam []	0.632	<b>0.582</b>	0.587	0.587	0.650
AMR Exam []	<b>0.692</b>	0.748	0.718	<b>0.682</b>	0.736

(b)

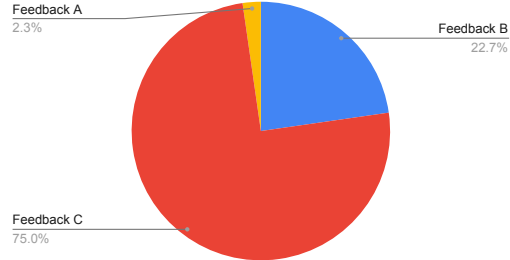
**Table 6:** Comparison of RMSE score between Random Forest (a) and AdaBoost (b) classifiers with language models (LM).

# Results

## Feedback: Survey Results

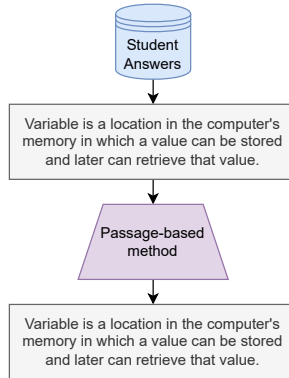
Question	What is a variable?
Reference Answer	A location in memory that can store a value.
Student Answer	a value/word that can assume any of a set of values
Feedback A	correct
Feedback B	missing keywords: location in memory
Feedback C	A variable is a location in memory that stores a value

Count of which feedback would one assign to this



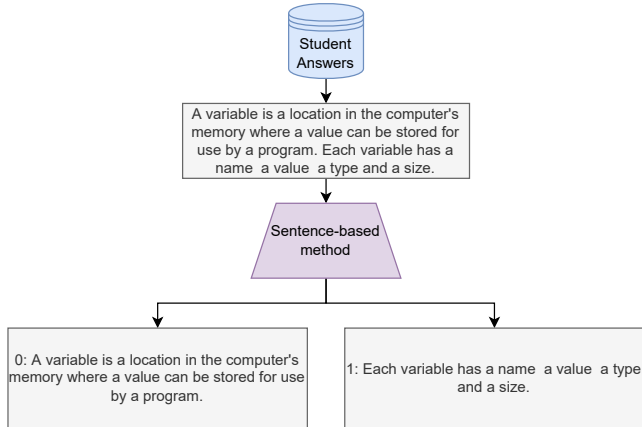
# Approach

## *Feature Extraction: Passage-based method*



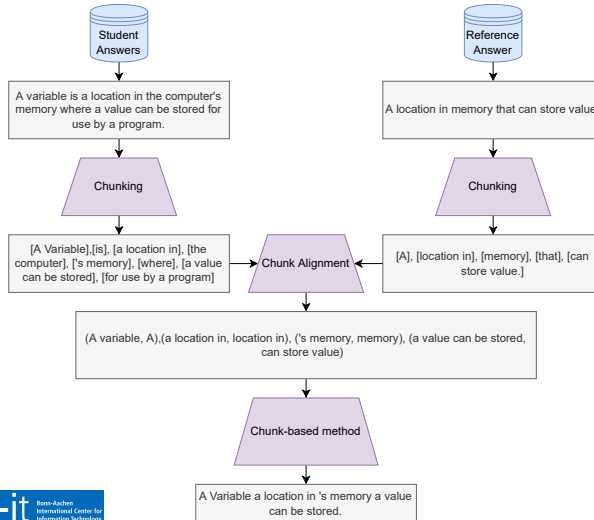
# Approach

## *Feature Extraction: Sentence-based method*



# Approach

## Feature Extraction: Chunk-based method





# Approach

## *Feature Extraction: RDF-based method*

