



Hochschule  
Bonn-Rhein-Sieg  
University of Applied Sciences



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

September 28, 2022

Elanton Fernandes

*Advisors*

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

# Agenda

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



# 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

- 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.

# Related Work

- Wu et al. (2021) designed a system to assign feedback called ProtoTransformer for evaluating programming based questions but not for short text answers. It used limited number of examples <sup>1</sup>.
- Ghavidel et al. (2020) passed raw text through a transformer as input and used the output of classification model (CLS) token as feature <sup>2</sup>.
- Mieskes and Pado, (2018) compared score assignment between automated and human assignment for RF, SVM, and DT classifiers across multiple datasets <sup>3</sup>.

---

<sup>1</sup> M. Wu, N. Goodman, C. Piech, and C. Finn, "Prototransformer: A meta-learning approach to providing student feedback," 2021.

<sup>2</sup> H. Ghavidel, A. Zouaq, and M. Desmarais, "Using bert and xlnet for the automatic short answer grading task," in Proceedings of the 12th International Conference on Computer Supported Education - Volume 1: CSEDU, pp. 58–67, INSTICC, SciTePress, 2020.

<sup>3</sup> M. Mieskes and U. Padó, "Work smart - reducing effort in short-answer grading," in Proceedings of the 7th workshop on NLP for Computer Assisted Language Learning, (Stockholm, Sweden), pp. 57–68, LiU Electronic Press, Nov. 2018.

# Dataset

Dataset	Domain	No. of Question Pairs	No. of Responses
Mohler <sup>4</sup>	Computer Science	81	2237
NN Exam <sup>5</sup>	Neural Network & AI	40	1137
AMR Exam <sup>6</sup>	Robotics	5	190

Table 1: Datasets used in score and feedback evaluation.

---

<sup>4</sup> M. Mohler, R. Bunesu, and R. Mihalcea, "Learning to grade short answer questions using semantic similarity measures and dependency graph alignments," June 2011.

<sup>5</sup> P. G. Plöger, "Neural network exam dataset," 2020.

<sup>6</sup> N. Hochgeschwender, "Autonomous mobile robots exam dataset," 2021.

# Approach

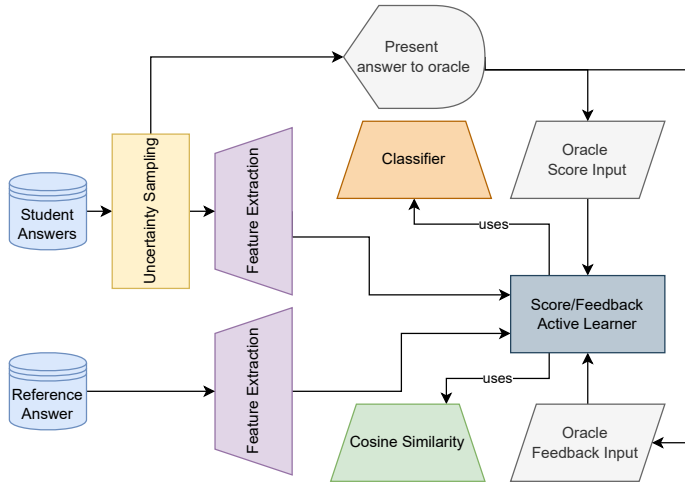
## *Contributions*

- Implement four methods to alter text for feature extraction.
- Implement feedback assignment for short text answers.
- Compare performance with five pre-trained language models and two classifiers.



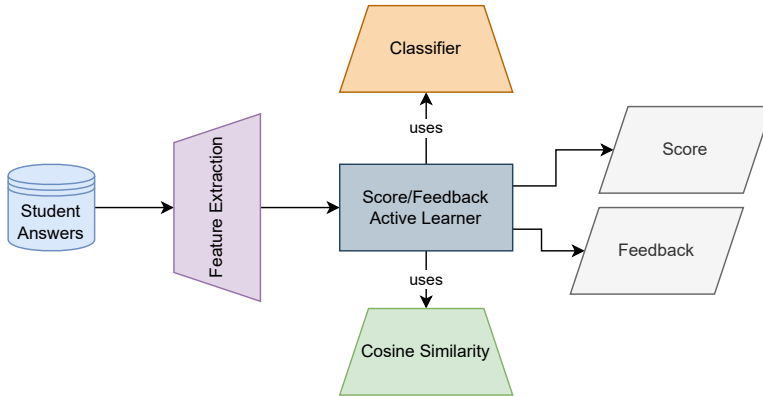
# Approach

## Training Cycle



# Approach

## *Prediction Cycle*



# Approach

## Uncertainty Sampling

Uncertainty sampling <sup>7</sup> 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:

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

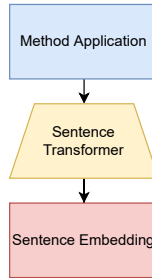
Where  $x$  is the feature,  $y$  is the class label prediction, and  $\hat{y} = \operatorname{argmax}_y P(y|x; \theta)$  is the class label that has the largest posterior probability using model  $\theta$ .

---

<sup>7</sup> B. Settles, "Computer Sciences Department Active Learning Literature Survey," 2009.

# Approach

## *Feature Extraction: Overview*



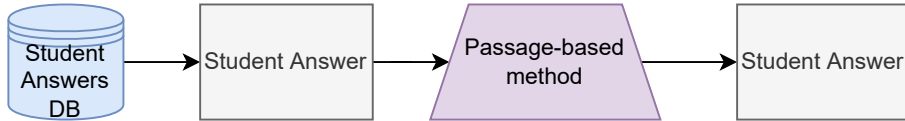
8

---

<sup>8</sup> N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," 2019.

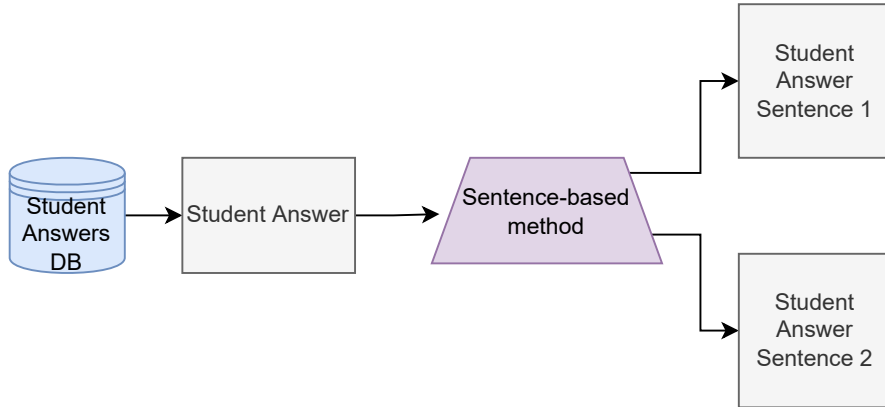
# Approach

## *Feature Extraction: Passage-Based Method*



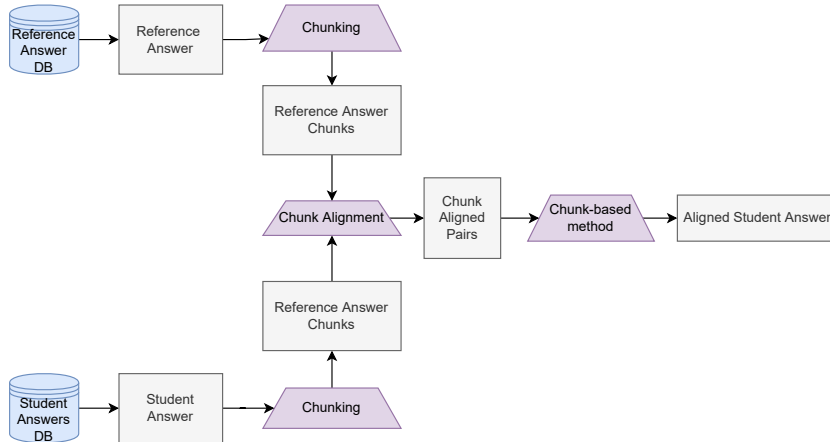
# Approach

## *Feature Extraction: Sentence-Based Method*



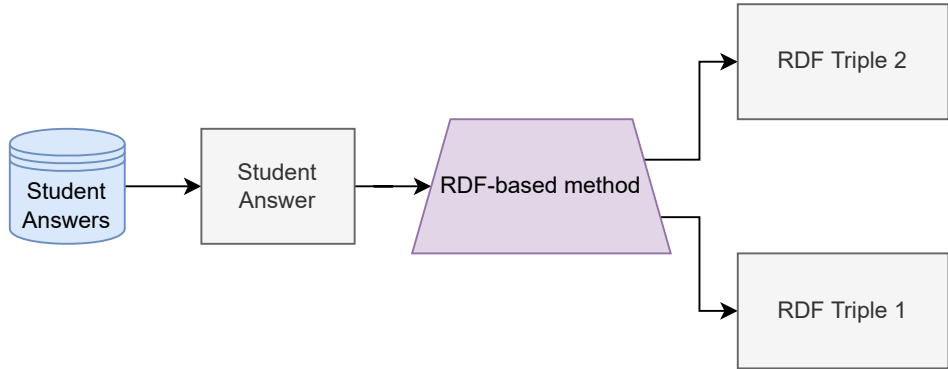
# Approach

## Feature Extraction: Chunk-Based Method



# Approach

## *Feature Extraction: Resource Description Framework (RDF) Based Method*





# Approach

## Language Models

Language Model <sup>9</sup>	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 2:** Displays pre-trained language models with their base model used in training and number of training tuples used.

<sup>9</sup> N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," 2019.

# Evaluation

## Score

- Pearsons Correlation

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

- RMSE Score

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

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

# 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 3: Presented survey to participants.

$$\text{Agreement Score} = \frac{\text{Model generated most rated feedback}}{\text{Total Number of Participants}}$$

# Results

## Notations

Method	Notation
Passage-based Methods	M1
Sentence-based Method	M2
Chunk-based Method	M3
RDF-based Method	M4

Language Model	Notation
all-mpnet-base-v2	LM1
all-distilroberta-v1	LM2
all-MiniLM-L12-v2	LM3
multi-qa-distilbert-cos-v1	LM4
all-MiniLM-L6-v2	LM5

# Results

Score: *Pearson Correlation (Methods)*

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

(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	0.428

(b)

**Table 4:** 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><u>0.802</u></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><u>0.760</u></b>
AMR Exam	0.453	<b>0.518</b>	<b><u>0.525</u></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 5:** 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><u>0.893</u></b>	<b>0.949</b>	<b>0.920</b>	0.942
NN Exam	<b><u>0.296</u></b>	<b>0.520</b>	<b>0.433</b>	<b>0.522</b>
AMR Exam	<b><u>0.596</u></b>	0.716	<b><u>0.596</u></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 6: 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><u>0.931</u></b>	<b>0.941</b>	<b>0.941</b>	<b>0.941</b>	<b>0.956</b>
NN Exam	<b><u>0.484</u></b>	0.591	<b>0.558</b>	<b>0.490</b>	<b>0.492</b>
AMR Exam	0.735	<b>0.680</b>	<b><u>0.676</u></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 7:** Comparison of RMSE score between Random Forest (a) and AdaBoost (b) classifiers with language models (LM).

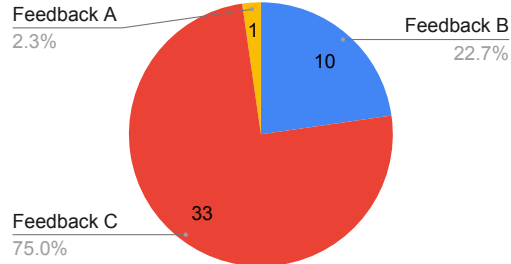


# Results

## 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 answer?



# Results

## *Feedback: Agreement Scores (Methods)*

Classifier	Methods			
	M1	M2	M3	M4
Random Forest	<b>60.00</b>	22.73	31.82	35.91
AdaBoost	<b>60.00</b>	22.73	31.82	35.91

Table 8: Mean agreement scores for Random Forest (a) and AdaBoost classifier (b) with methods.

# Results

## *Feedback: Agreement Scores (Models)*

Classifier	Language Models				
	LM1	LM2	LM3	LM4	LM5
Random Forest	25.11	26.82	24.66	<b>37.05</b>	21.25
AdaBoost	25.11	26.82	24.66	<b>37.05</b>	21.25

Table 9: Mean agreement scores for Random Forest and AdaBoost classifier with Language Models.

# Results

## Summary: Scores

Dataset	Method	Model	Classifier
Mohler	M1	LM1	Random Forest
NN Exam	M1	LM5	Random Forest
AMR Exam	M1	LM3	Random Forest

(a)

Dataset	Method	Model	Classifier
Mohler	M1	LM1	Random Forest
NN Exam	M1	LM1	Random Forest
AMR Exam	M1& M3	LM3	Random Forest

(b)

Table 10: Performance summary Pearson correlation (a) and RMSE score (b)

# Results

## *Feedback*

Dataset	Method	Model	Method-Model	Classifier
Mohler	M1	LM4	M1-LM4	Random Forest

Table 11: Results of feedback evaluation

# Summary

- Four methods implemented to alter student answer text.
- Designed system that keeps oracle in training loop.
- System allows assignment of feedback.
- Pearson correlation and RMSE score used as score evaluation metrics.
- Survey created and used in the evaluation of the feedback assigned by the model.

# Future Work

- RDF-based method did not give good RDF triplets. Fine tune model for better results.
- Fine tune language models similar to [multi-qa-distilbert-cos-v1](#) to create embeddings.
- Chunk-based Method took 41 sec per prediction. This time needs to be reduced.





# Foundation

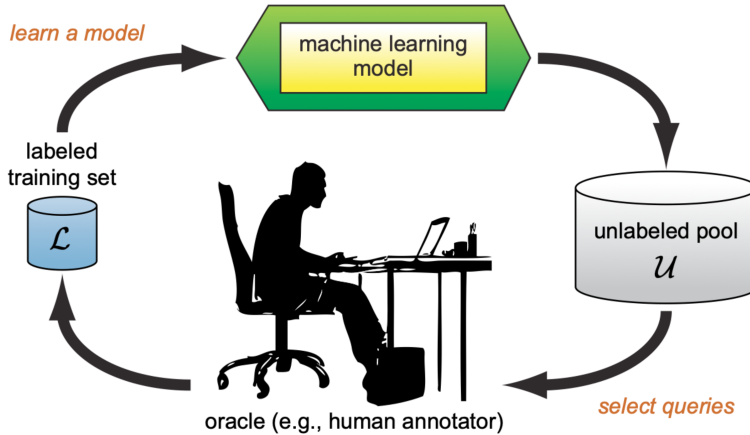
## *Cosine Similarity*

If  $x$  and  $y$  are two sentences and  $e_x$  and  $e_y$  are embeddings of these sentences. Then the cosine similarity is given by:

$$S_{\text{cosine}}(e_x, e_y) = \frac{e_x \cdot e_y}{\|e_x\| \|e_y\|} \quad (4)$$

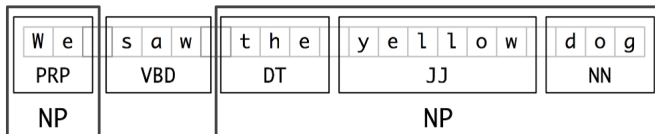
# Foundation

## Active Learning



# Foundation

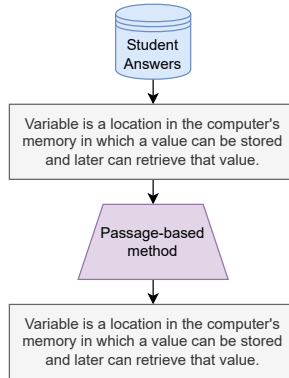
## Chunking



Symbol	Meaning	Symbol	Meaning
NP	Noun Phrase	VP	Verb phrase
NN	Noun	DT	zero or one determiner
JJ	One or more adjectives	PRP	Preposition

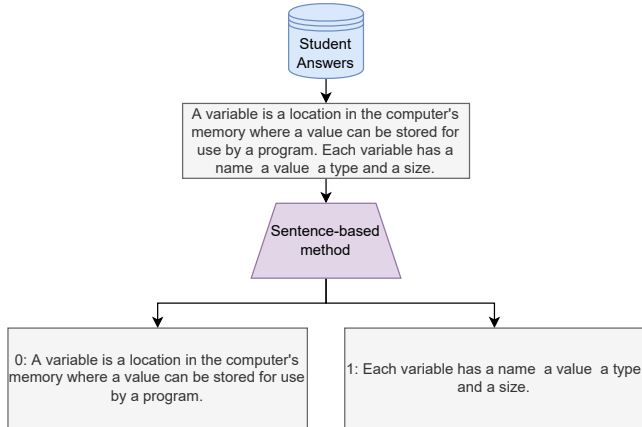
# Approach:Example

## *Feature Extraction: Passage-based method*



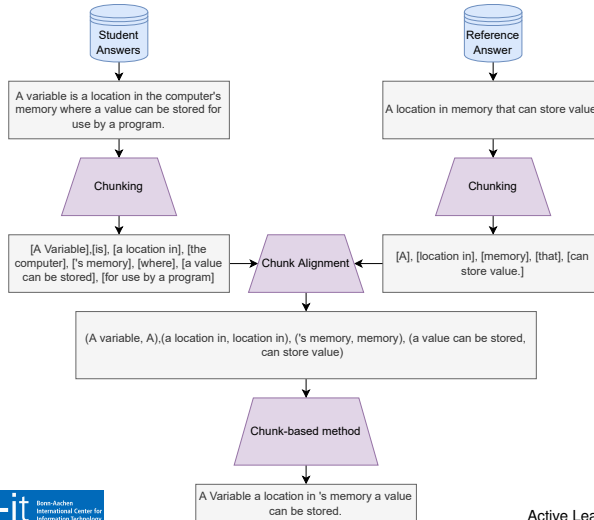
# Approach

## *Feature Extraction: Sentence-based method*



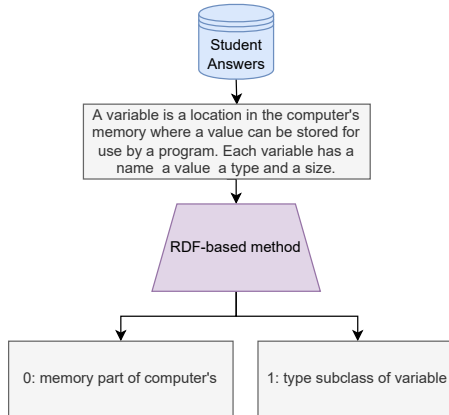
# Approach

## Feature Extraction: Chunk-based method



# Approach

## *Feature Extraction: RDF-based method*



# Results

## *Feedback: Survey Results*

Question	What stages of the software lifecycle are influenced by the testing stage
Reference Answer	The testing stage can influence both the coding stage (phase 5) and the solution refinement stage (phase 7)
Student Answer	All stages are influenced except setting the program requirements. If a test fails it can change the whole design implementation etc of a program as well as the final outcome.
Feedback A	The testing phase affects the coding/production phase and the refinement/maintenance phase
Feedback B	Correct
Feedback C	Missing keywords: Refinement stage,Coding phase



Count of Which feedback would you assign to this answer?

Feedback B

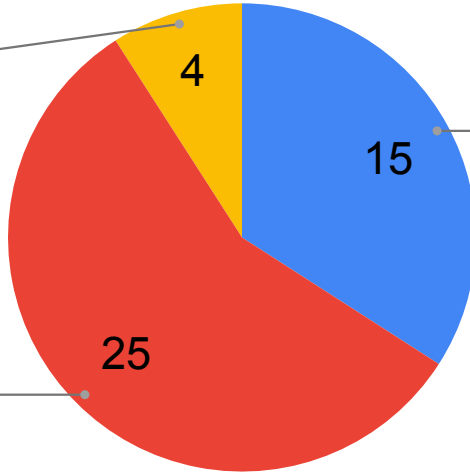
9.1%

Feedback A

56.8%

Feedback C

34.1%



# Results

## *Feedback: Survey Results*

Question	What are the main advantages associated with object oriented programming?
Reference Answer	Abstraction and reusability.
Student Answer	Re-usability and ease of maintenance
Feedback A	Missing keywords: Reusability, Abstraction
Feedback B	Correct
Feedback C	The main advantages of OOP are abstraction and reusability
Feedback D	Encapsulation is similar to abstraction but not the same and second advantage is reusability.

Count of Which feedback would you assign to this answer?

Feedback A

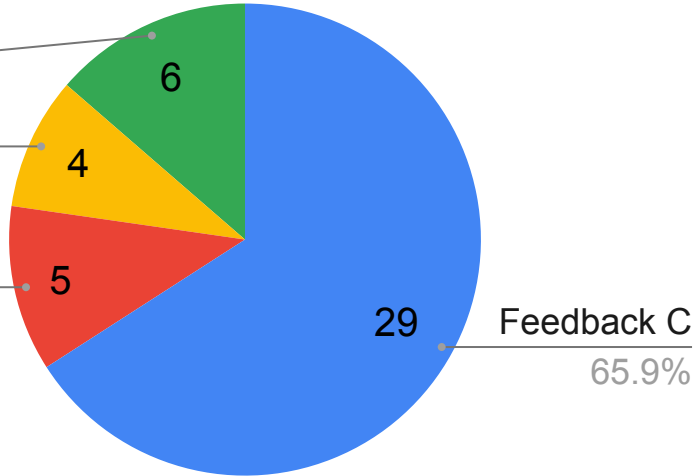
13.6%

Feedback B

9.1%

Feedback D

11.4%



# 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 answer?

Feedback A

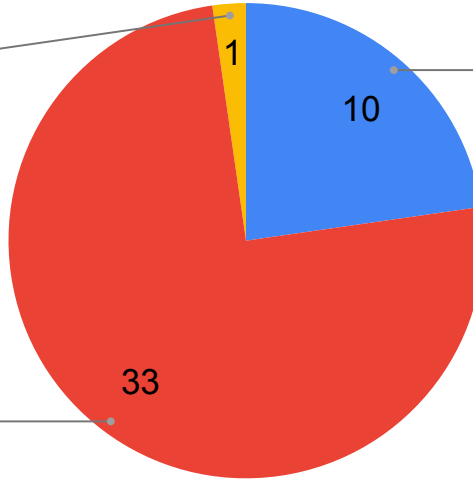
2.3%

Feedback B

22.7%

Feedback C

75.0%



# Results

## *Feedback: Survey Results*

Question	What is the difference between a constructor and a function?
Reference Answer	A constructor is called whenever an object is created whereas a function needs to be called explicitly. Constructors do not have return type but functions have to indicate a return type.
Student Answer	Constructors don't have a return type.
Feedback A	Missing point: A constructor is called whenever an object is created whereas a function needs to be called explicitly.
Feedback B	Correct
Feedback C	Missing point: Constructors do not have return type but functions have to indicate a return type.
Feedback D	Answer not explained properly. Information on function/ constructor calling and their return type missing.

Count of Which feedback would you assign to this answer?

Feedback C

9.1%

Feedback B

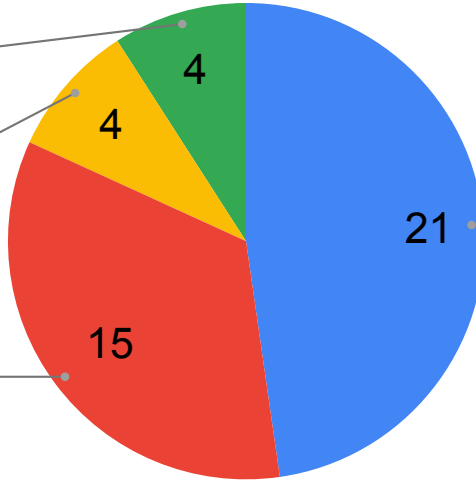
9.1%

Feedback D

34.1%

Feedback A

47.7%



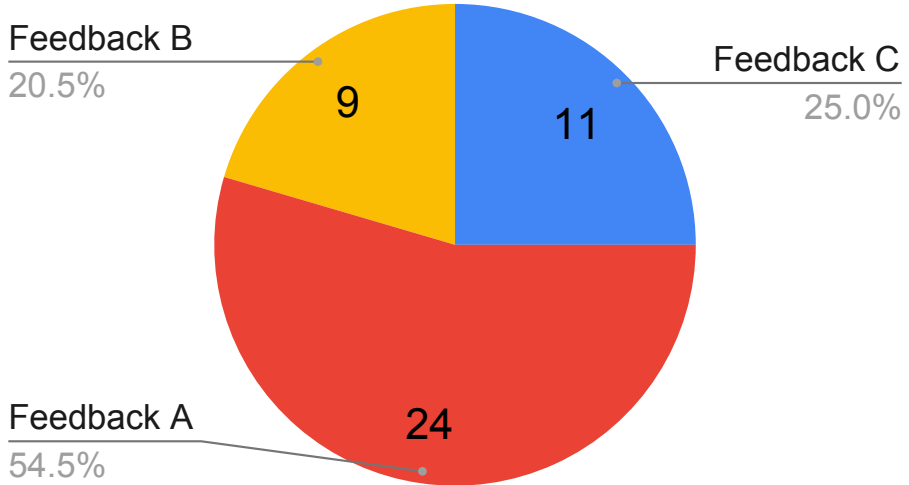
# Results

## *Feedback: Survey Results*

Question	When does C++ create a default constructor?
Reference Answer	If no constructor is provided the compiler provides one by default. If a constructor is defined for a class the compiler does not create a default constructor.
Student Answer	When non are provided
Feedback A	If no constructor is defined then a default constructor is created during compilation.
Feedback B	Missing keywords: during compilation
Feedback C	Correct



Count of Which feedback would you assign to this answer?



# Results

## *Prediction Time*

Method	Avg. Prediction Time (sec)
Passage-Based	<b>0.0344</b>
Sentence-Based	0.0749
Chunk-Based	41.3039
RDF-Based	5.8424

**Table 12:** Shows average prediction time (in seconds) per response for each method.