

ReflectSumm: A Benchmark for Course Reflection Summarization

Yang Zhong*, Mohamed Elaraby*, Diane Litman
Ahmed Ashraf Butt, Muhsin Menekse

*Equal Contribution



Text Summarization

ML models performed well on standard benchmarks with ample data



Text Summarization

ML models performed well on standard benchmarks with ample data



Yet, little is known about how models work on low-resource real-life applications



I found sorting algorithm most confusing because .

Class goes too fast ...

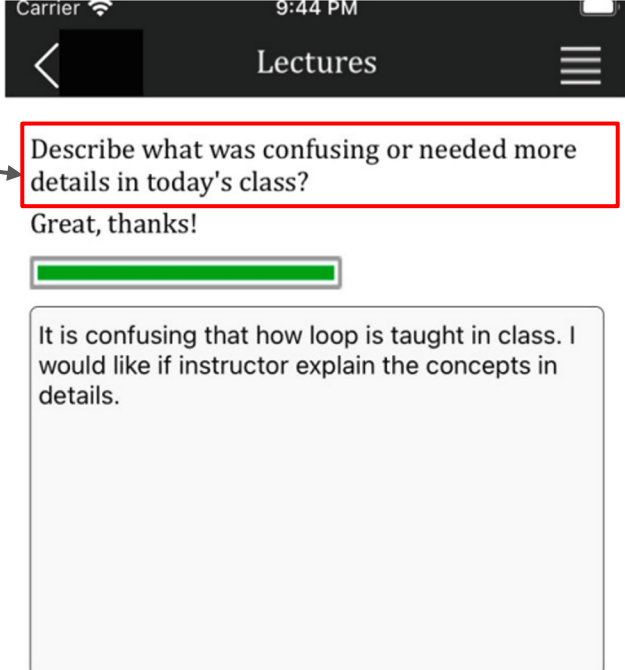
... Analysis of Time Complexity ...

Summary of Reflections



CourseMIRROR Data Collection

- **Step 1:** Students are prompted to write a reflection about what they found confusing or interesting



Carrier 9:44 PM

< Lectures ≡

Describe what was confusing or needed more details in today's class?

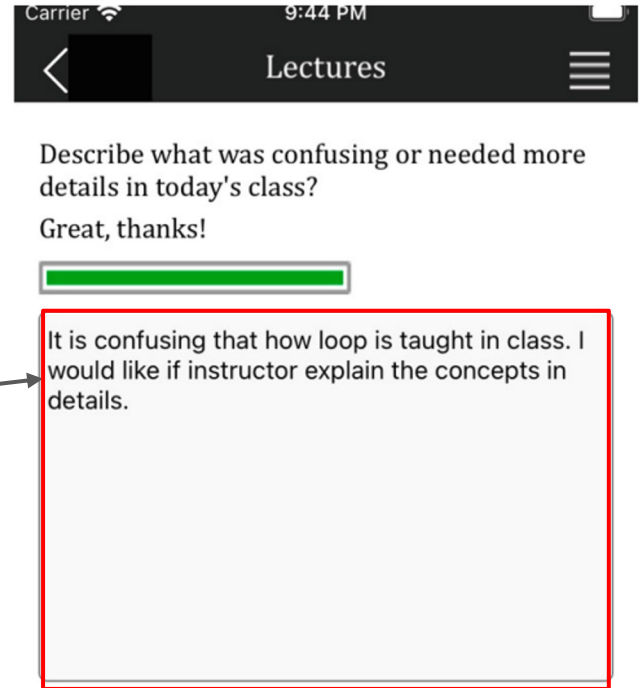
Great, thanks!

Progress bar (green)

It is confusing that how loop is taught in class. I would like if instructor explain the concepts in details.

CourseMIRROR Data Collection

- Step 1: Students are prompted to write a reflection about what they found confusing or interesting
- **Step 2:** Reflections are collected through the CourseMIRROR app



Carrier 9:44 PM

< Lectures ≡

Describe what was confusing or needed more details in today's class?
Great, thanks!

It is confusing that how loop is taught in class. I would like if instructor explain the concepts in details.

The screenshot shows a mobile app interface for 'Lectures'. At the top, there's a status bar with 'Carrier' and '9:44 PM'. Below that is a navigation bar with a back arrow, the title 'Lectures', and a hamburger menu icon. The main content area contains a text prompt: 'Describe what was confusing or needed more details in today's class?' followed by 'Great, thanks!'. Below the prompt is a green progress bar. A red box highlights a text input field containing the student's reflection: 'It is confusing that how loop is taught in class. I would like if instructor explain the concepts in details.' An arrow points from the 'Step 2' bullet point in the text to this red box.

Course Reflection Summarization

1. Offering reflection summaries can aid instructors in **enhancing lecture preparation and supporting students**

Reflection Prompt

Describe what you found most interesting in today's class

Student Reflections

- Nothing in particular today
- Despite the confusion, I did find setting up these problems to be very interesting and rewarding.
- Equipotentials
- i thought the breakout room questions were interesting because i learned how to do questions
- I found the last problem in class the most interesting because it was proven we can derive almost anything
- The most interesting thing was that finding electric potential doesn't require a path, but only the magnitude of the charge and it's distance from the point of interest.
- I really enjoy line integrals and I can tell that we're moving towards using them to calculate potential.
- Collection of point charges (pairing them)
- How we can calculate something so complicated as electrons passing through an area is very cool.
- I found equipotentials to be the most interesting thing, especially drawing a equipotentials for a dipole!
- I thought it was interesting that V_{net} is equal to all V_s added together
- I found how conductors act to be interesting.

Course Reflection Summarization

1. Offering reflection summaries can aid instructors in enhancing lecture preparation and supporting students

2. Course reflections **vary in the length and structure** (different from opinion summarization)

Reflection Prompt

Describe what you found most interesting in today's class

Student Reflections

- Nothing in particular today
- Despite the confusion, I did find setting up these problems to be very interesting and rewarding.
- Equipotentials
- I thought the breakout room questions were interesting because I learned how to do questions
- I found the last problem in class the most interesting because it was proven we can derive almost anything
- The most interesting thing was that finding electric potential doesn't require a path, but only the magnitude of the charge and its distance from the point of interest.
- I really enjoy line integrals and I can tell that we're moving towards using them to calculate potential.
- Collection of point charges (pairing them)
- How we can calculate something so complicated as electrons passing through an area is very cool.
- I found equipotentials to be the most interesting thing, especially drawing a equipotentials for a dipole!
- I thought it was interesting that V_{net} is equal to all V_s added together
- I found how conductors act to be interesting.

Course Reflection Summarization

1. Offering reflection summaries can aid instructors in enhancing lecture preparation and supporting students
2. Course reflections vary in the length and structure (different from opinion summarization)
3. **Automatic summarization** can help **scale the use** of reflections in educational practice

Lack of Benchmarks!

Reflection Prompt

Describe what you found most interesting in today's class

Student Reflections

- Nothing in particular today
- Despite the confusion, I did find setting up these problems to be very interesting and rewarding.
- Equipotentials
- i thought the breakout room questions were interesting because i learned how to do questions
- I found the last problem in class the most interesting because it was proven we can derive almost anything
- The most interesting thing was that finding electric potential doesn't require a path, but only the magnitude of the charge and it's distance from the point of interest.
- I really enjoy line integrals and I can tell that we're moving towards using them to calculate potential.
- Collection of point charges (pairing them)
- How we can calculate something so complicated as electrons passing through an area is very cool.
- I found equipotentials to be the most interesting thing, especially drawing a equipotentials for a dipole!
- I thought it was interesting that V_{net} is equal to all Vs added together
- I found how conductors act to be interesting.

Our Work

- Manually annotated corpus with different types of summaries and rich metadata
- A comprehensive benchmarking of models for varied summarization tasks

Our Work

- **Manually annotated corpus with different types of summaries and rich metadata**
- A comprehensive benchmarking of models for varied summarization tasks

ReflectSumm Dataset

Metadata
Gender;
Ethnicity; Age

Reflection Prompt

Describe what you found most interesting in today's class

Student Reflections

- Nothing in particular today -> 1.0
- Despite the confusion, I did find setting up these problems to be very interesting and rewarding. -> 3.0
- Equipotentials -> 2.0
- i thought the breakout room questions were interesting because i learned how to do questions -> 4.0
- I found the last problem in class the most interesting because it was proven we can derive almost anything. -> 4.0
- ★ The most interesting thing was that finding electric potential doesn't require a path, but only the magnitude of the charge and it's distance from the point of interest. -> 4.0
- I really enjoy line integrals and I can tell that we're moving towards using them to calculate potential. -> 4.0
- Collection of point charges (pairing them) -> 2.0
- How we can calculate something so complicated as electrons passing through an area is very cool. -> 3.0
- ★ I found equipotentials to be the most interesting thing, especially drawing a equipotentials for a dipole! -> 4.0
- I thought it was interesting that V_{net} is equal to all Vs added together -> 4.0
- I found how conductors act to be interesting. -> 3.0

★ Extractive Summary

- I found equipotentials to be the most interesting thing, especially ...
- The most interesting thing was that finding electric potential doesn't require a path ...
(three more reflections selected from input)

Annotated
Specificity Score

ReflectSumm Dataset

Metadata
Gender;
Ethnicity; Age

Reflection Prompt

Describe what you found most interesting in today's class

Student Reflections

- Nothing in particular today -> 1.0
- Despite the confusion, I did find setting up these problems to be very interesting and rewarding. -> 3.0
- Equipotentials -> 2.0
- i thought the breakout room questions were interesting because i learned how to do questions -> 4.0
- I found the last problem in class the most interesting because it was proven we can derive almost anything. -> 4.0
- The most interesting thing was that finding electric potential doesn't require a path, but only the magnitude of the charge and it's distance from the point of interest. -> 4.0
- I really enjoy line integrals and I can tell that we're moving towards using them to calculate potential. -> 4.0
- Collection of point charges (pairing them) -> 2.0
- How we can calculate something so complicated as electrons passing through an area is very cool. -> 3.0
- I found equipotentials to be the most interesting thing, especially drawing a equipotentials for a dipole! -> 4.0
- I thought it was interesting that V_{net} is equal to all Vs added together -> 4.0
- I found how conductors act to be interesting. -> 3.0

Annotated
Specificity Score

Extractive
Summary

Abstractive
Summary

Phrase
Summary

- I found equipotentials to be the most interesting thing, especially ...
- The most interesting thing was that finding electric potential doesn't require a path ...
(three more reflections selected from input)

The students today found calculations and relationships to other concepts that they have learned in this and other classes interesting. They also found potential energy and equipotentials very interesting, as well as some integration concepts.

- equipotentials
- calculations
- relations to old concepts
- potential
- integration

ReflectSumm Dataset

- Student reflections on lectures collected from real-world lectures
- **782** input-summary pairs from 24 university courses, spanning four STEM subjects

*We only included data from students who explicitly provided consents, and carefully removed all private information

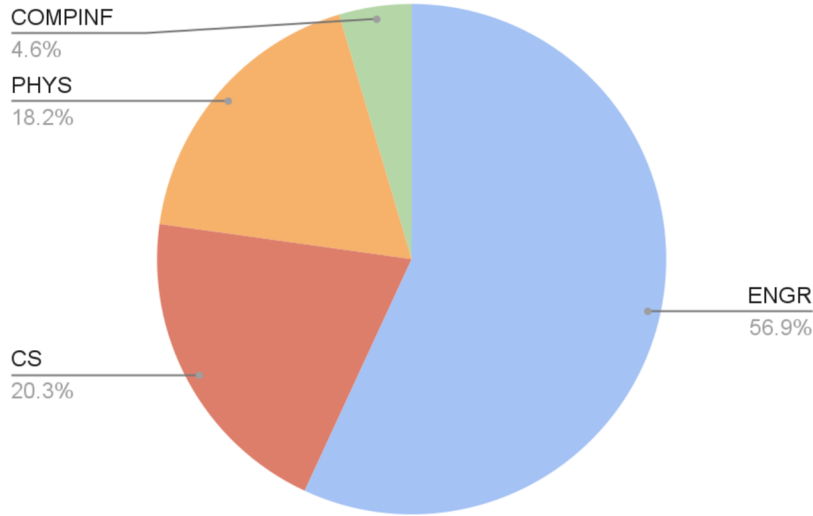
ReflectSumm Dataset

- Student reflections on lectures collected from real-world lectures
- **782** input-summary pairs from 24 university courses, spanning four STEM subjects
- **Manual annotations** for each pair of **input-summary**
 - **Individual input reflection** has annotated specificity score and metadata of the author (demographic information*)
 - **Three types of summaries** are manually annotated by high-quality in-house annotators.

*We only included data from students who explicitly provided consents, and carefully removed all private information

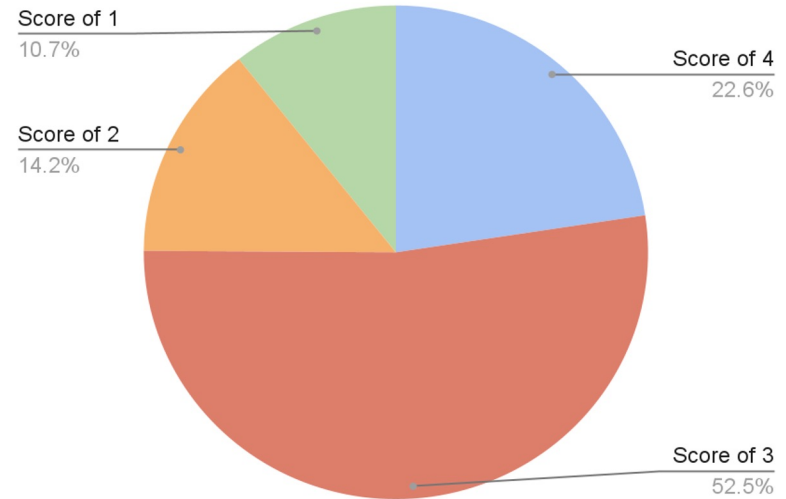
Dataset Analysis

Lecture Distribution



Our dataset covers diverse subjects

Reflection Specificity Distribution



Students write high quality reflections with details

Our Work

- Manually annotated corpus with different types of summaries and rich metadata
- **A comprehensive benchmarking of models for varied summarization tasks**

Extractive

Phrase

Abstractive

Extractive Summarization

Task: Select 5 reflections as the summary

Extractive

Phrase

Abstractive

Extractive Summarization

Task: Select 5 reflections as the summary

Models:

- **BERTSUM-EXT** (Liu and Lapta, 2019)
 - Fine-tuned on CNN/DM (287K) or ReflectSumm (less than 800)
- **MatchSUM** (Zhong et al., 2020)
- **ChatGPT (GPT)** with different prompts
 - Reflections only
 - Reflections + Specificity

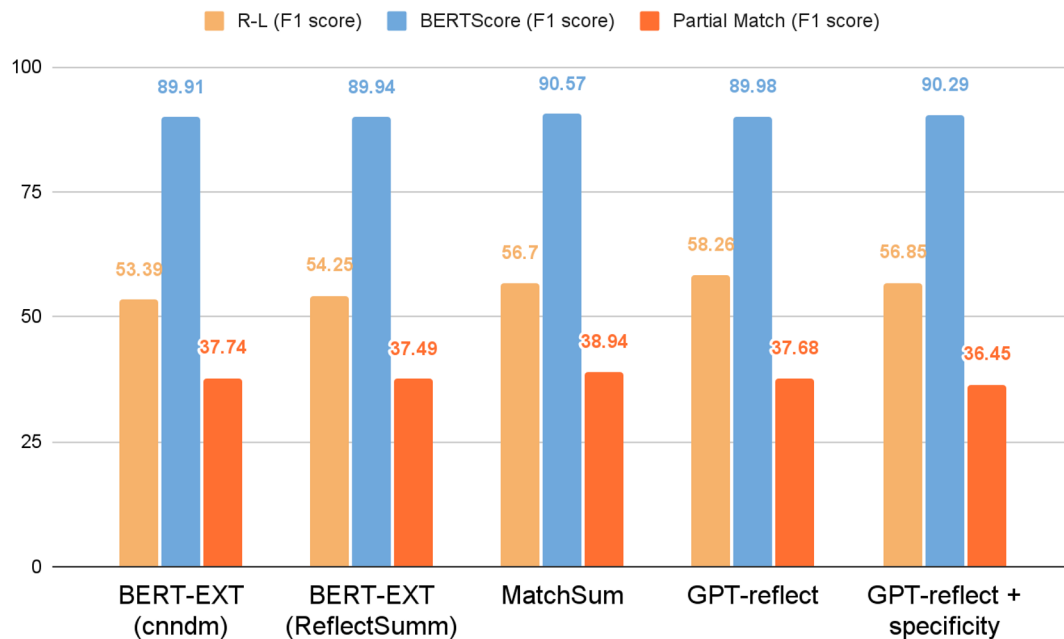
Eval Metrics: ROUGE (R-1/R-2/R-L); BERTScore; Partial Match F1

Extractive

Phrase

Abstractive

Extractive Summarization



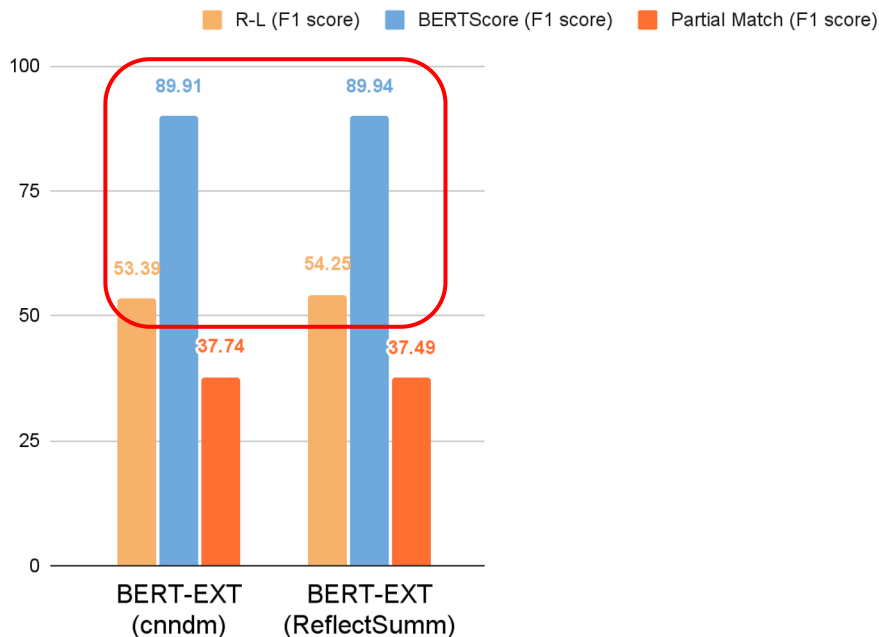
Extractive

Phrase

Abstractive

Extractive Summarization

The performance slightly improves when utilizing the BERT-EXT model trained on our dataset, which is much smaller than the CNN/DM dataset



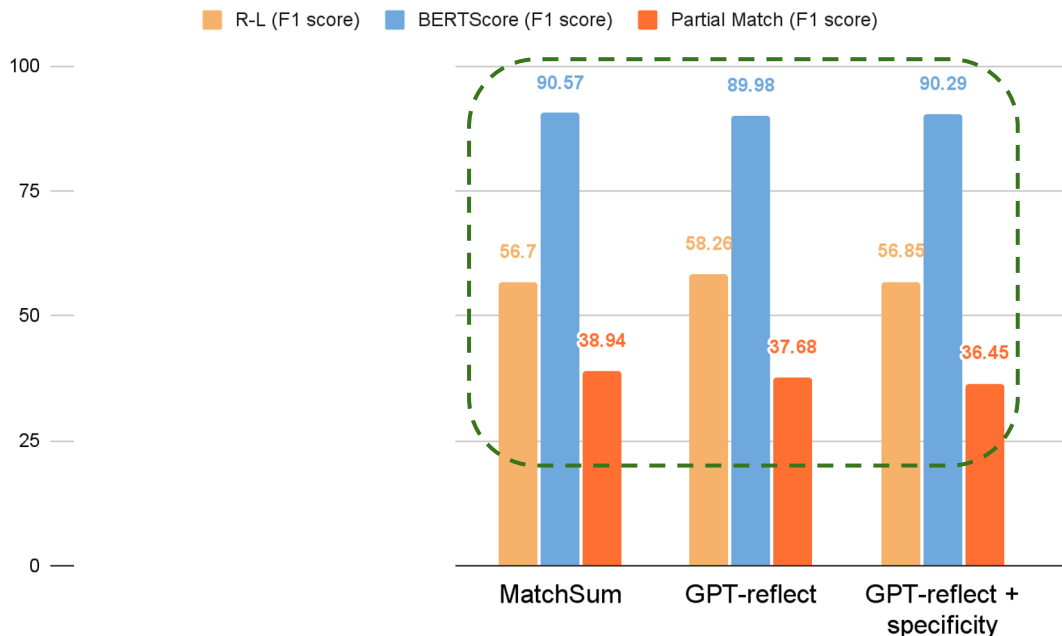
Extractive

Phrase

Abstractive

Extractive Summarization

The zero-shot GPT-based model achieves comparable or superior performance compared to robust baselines, but it may occasionally struggle to faithfully extract complete reflections



Extractive

Phrase

Abstractive

Phrase Summarization

Task: Generate 5 phrases to summarize the reflections

Models:

- **PhraseSum:** deployed baseline (Luo and Litman 2015)
- **ChatGPT**
 - GPT-noun: prompt to generate noun phrases
 - GPT-Human: A more intricate prompt with human guidelines

Eval Metrics: ROUGE (R-1/R-2/R-L); BERTScore

Extractive

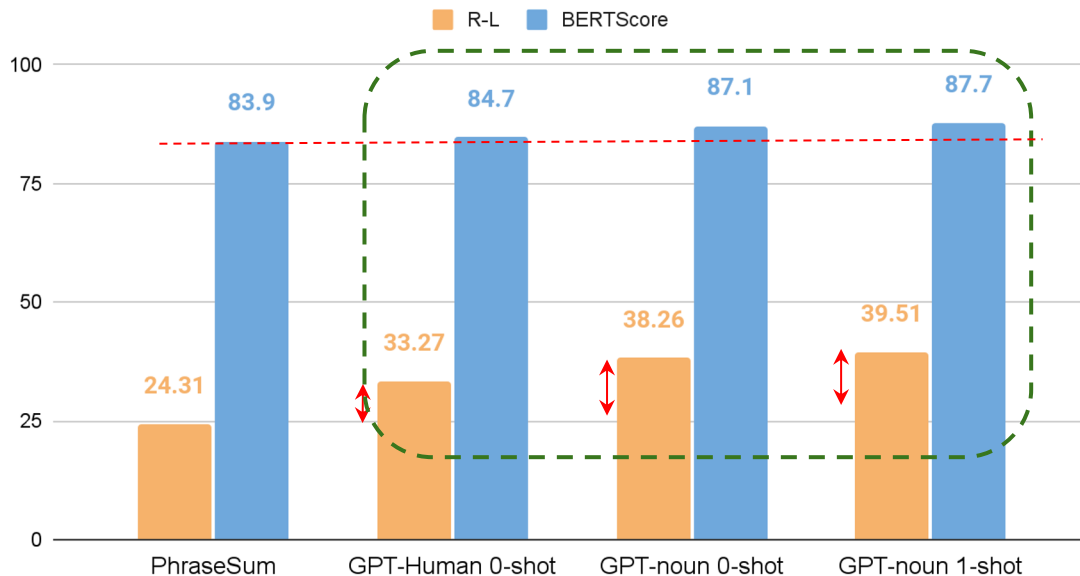
Phrase

Abstractive

Phrase Summarization

GPT models outperform the baseline by large margin

GPT-based models are sensitive to the prompt, and few-shot could help produce higher-quality phrases



Extractive

Phrase

Abstractive

Abstractive Summarization

Finetune pretrained encoder-decoder



- Vanilla fine-tuning
- Fine-tuning with reflection specificity

Prompt GPT-3.5-turbo



- Zeroshot prompting
- Zeroshot prompting reflection specificity

Extractive

Phrase

Abstractive

Abstractive Summarization

Evaluation

- **Reference-based metrics**
ROUGE (R-1/R-2/R-L); and BERTScore
- **Factuality metrics**
Entailment based approaches — SummaC (Laban et al., 2022)

Extractive

Phrase

Abstractive

Abstractive Summarization

Reference-based metrics

Specificity information showed positive influence in case of fine-tune pretrained models

| Model | R-1 | R-2 | R-L | BS |
|---------------|--------------|--------------|---------------|--------------|
| BART-Large | 47.09 | 24.17 | 43.76 | 90.49 |
| + specificity | 47.70 | 24.85 | 44.41* | 90.57 |
| GPT-Human | 35.83 | 9.40 | 31.85 | 88.23 |
| + specificity | 36.73 | 9.13 | 31.64 | 88.27 |
| GPT-one-shot | 36.86 | 9.46 | 31.96 | 88.26 |

Extractive

Phrase

Abstractive

Abstractive Summarization

Factuality Metrics: SummaC

Is GPT more factual than human reference?

| Model | SUMMAC ↑ | |
|-----------------|-------------|-------------|
| | Sentence | Document |
| Human-reference | 0.25 | 0.22 |
| BART-Large | 0.25 | 0.21 |
| + specificity | 0.25 | 0.22 |
| GPT-Human | 0.26 | 0.31 |
| + specificity | 0.27 | 0.26 |
| GPT-one-shot | 0.26 | 0.26 |

Extractive

Phrase

Abstractive

Abstractive Summarization

Example

- One thing I found interesting was how many categories of machine learning there are.
- Supervised and unsupervised learning as well as discrete and continuous labels and how they all related to one another.
- Different categories of machine learning.
- The relationship between unsupervised and supervised deep learning.

Source Document

Extractive

Phrase

Students enjoyed learning about the differences between supervised and unsupervised learning. Along with that, they also enjoyed learning about the different categories in Machine Learning and the different categorization and classification methods.

Factual Summary

Abstractive

Abstractive Summarization

Example

**SummaC: 0.2
(low factuality)**

- One thing I found interesting was how many categories of machine learning there are.
- Supervised and unsupervised learning as well as discrete and continuous labels and how they all related to one another.
- Different categories of machine learning.
- The relationship between unsupervised and supervised deep learning.

Students enjoyed learning about the differences between supervised and unsupervised learning. Along with that, they also enjoyed learning about the different categories in Machine Learning and the different categorization and classification methods.

Source Document

Factual Summary

Extractive

Phrase

Abstractive

Abstractive Summarization

Example

SummaC: 0.2
(low factuality)

- One thing I found interesting was how many categories of machine learning there are.
- Supervised learning is used as well as unsupervised learning. How the relationship between the two is different.
- The relationship between unsupervised and supervised deep learning.

We need to develop **better factuality metrics** for these types of documents

Students are learning about the supervised and unsupervised learning categories in the different methods.

Source Document

Factual Summary

Extractive

Phrase

Abstractive

Broader Impact

- **Rich metadata** (demographic information of student) can be applied for studies on fairness / equity issues
- Our models can **enable new downstream functionalities** such as generating recommended readings and explaining confusing concepts based on summary output



For instance, our dataset can enable researchers to analyze **differences in reflection submission rates** and **specificity among different student groups**

Broader Impact

- Rich metadata (demographic information of student) can be applied for studies on fairness / equity issues.
- Our models can enable new downstream functionalities such as generating recommended readings and explaining confusing concepts based on summary output.



Future Work

- Investigating improved **prompting techniques**.
- Extending the corpus to **other subject domains** (i.e. Psychology).
- Designing a **dynamic system** capable of adjusting the quantity of extracted summary outputs according to the size of the lecture.

Conclusion

- Manually annotated corpus with **different types of summaries** and **rich metadata** (age, gender, ethnicity)

Extractive

Phrase

Abstractive

- A **comprehensive benchmarking** of models for varied summarization tasks

- Check the prompts, system outputs, and dataset at <https://github.com/EngSalem/ReflectSUMM>

- Contact:

- Yang Zhong yaz118@pitt.edu
- Mohamed Elaraby mse30@pitt.edu

