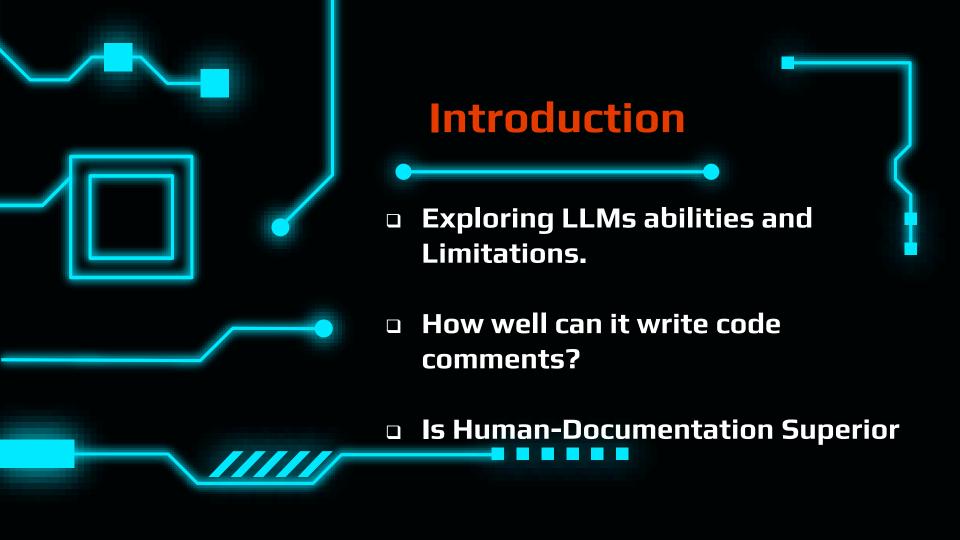
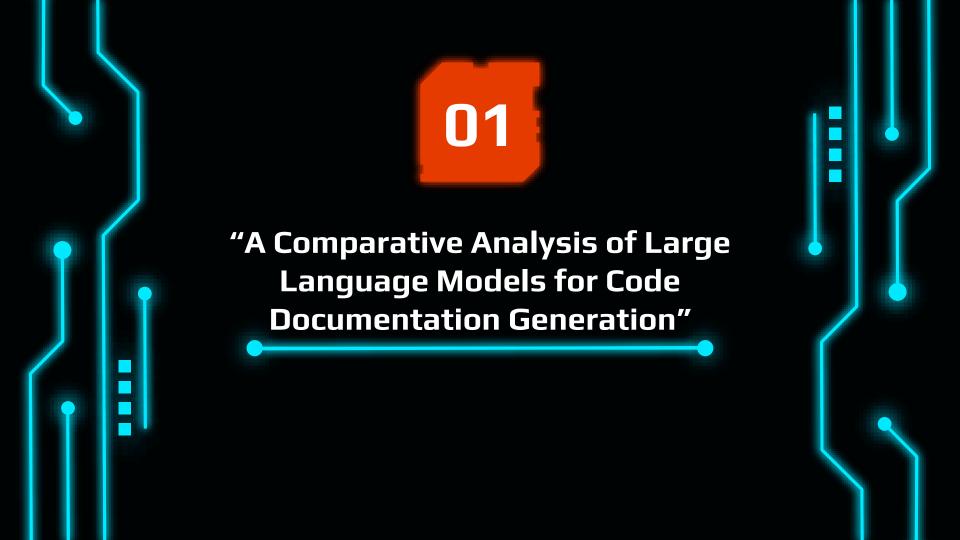


ChatGPT's Ability to generate comments for codes

Hassan Abdul-Mohsen Abdullah Mishary







Idea

- Evaluate the effectiveness of LLMs in generating code documentation
 - comparing Ai models across different documentation levels (inline, function, and file-level).

AI Models Used:





- □ checklist-based evaluation approach
- □ 14 Python code snippets from welldocumented open-source projects

(Bitcoin, TensorFlow, SQLAlchemy, Scikit-Learn, and Reddit)

- generate documentation using prompts
- two expert reviewers assess the output

Evaluation metrics used:

O1 Accuracy

04

Understandability

02 Completeness

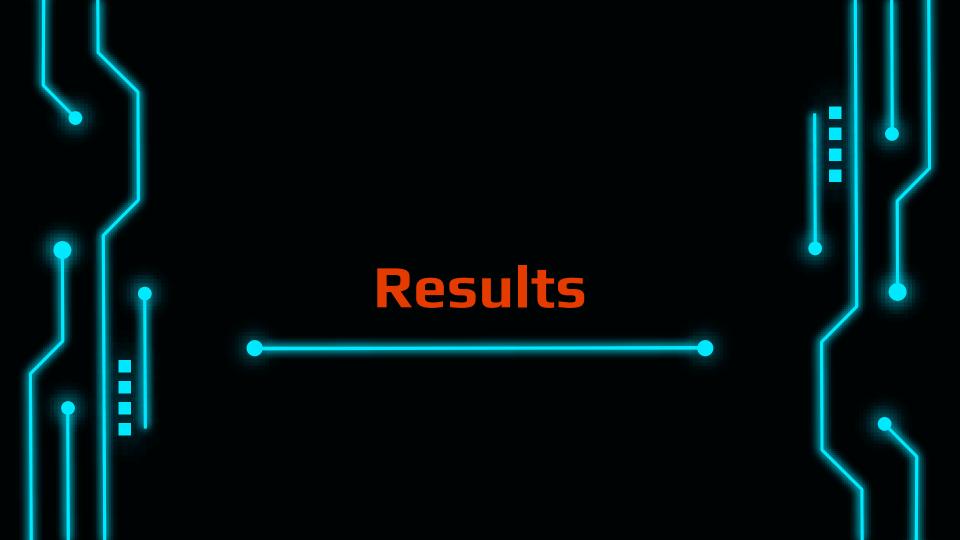
05

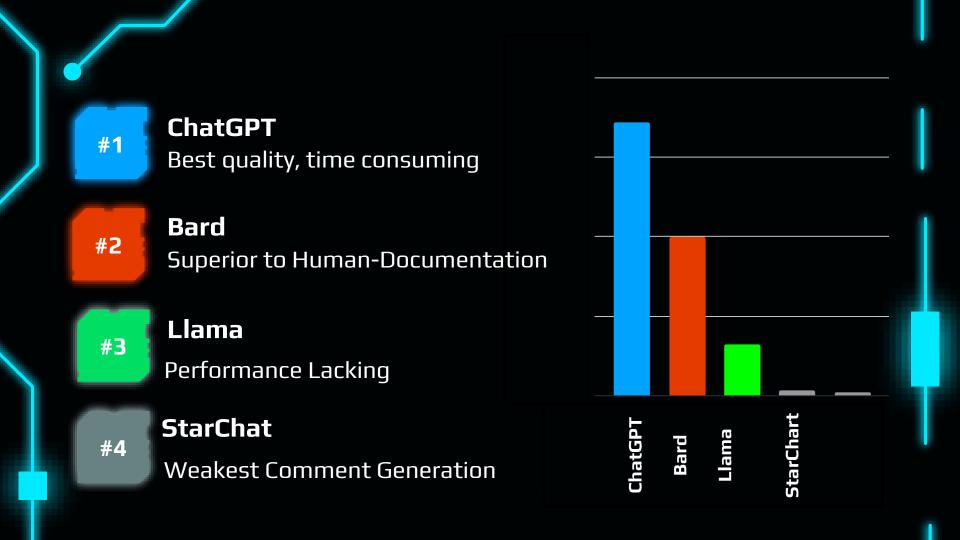
Readability

03 Relevance

06

Time Taken









Idea

- ☐ investigate ChatGPT's effectiveness in automatic code summarization
- comparing performance against state-ofthe-art (SOTA) models

- Design prompts to guide ChatGPT in generating comments for code snippets
- Utilize the CSN-Python dataset for code test set.
- Compare summaries with those produced by SOTA models

Evaluation Metrics

- measure the quality of generated text against reference summaries, using the following metrics:
- BLEU
- METEOR
- ROUGE-L

able 5: Performance of ChatGPT with different prompts

Method	CSN-Python		
	BLEU	METEOR	ROUGE-L
ChatGPT (short)	4.85	12.90	14.41
ChatGPT (one sentence)	10.66	13.36	20.37
ChatGPT (within 40 words)	4.83	13.01	14.05

e 6: Performance of ChatGPT and three SOTA models

Method	CSN-Python		
	BLEU	METEOR	ROUGE-L
NCS	15.8	10.6	31.3
CodeBERT	18.7	12.4	34.8
CodeT5	20.0	14.7	37.7
ChatGPT (one sentence)	10.28	14.40	20.81

Results

- ChatGPT's performance was significantly lower than SOTA models.
- ChatGPT's code summaries were less aligned with reference summaries than those produced by other models

```
# checking response.status_come (if you go 5
 if response.status_code != 200:
       print(f"Status: {response.status_code; - Try research
  else:
        print(f"Status: (response.status_code \ur)
 # using BeautifulSoup to parts the recomme
 SOUP = BeautifulSoup(response.content, "non-neural
images = soup.find_all("asg", attractates was seen
 # downloading images
```





Idea

- ☐ Investigates GPT-4's effectiveness
- generating Java class-level and methodlevel documentation
- compare to manually written reference documentation.

- using predefined prompts.
- GPT-4 generates Javadoc comments for class and method levels.
 - Output is evaluated using both automated (BLEU, ROUGE) and manual methods.
 - □ Select Eleven Java classes and 55 methods from open-source projects

Evaluation metrics used:

01 Correctness

04 BLEU

02 Completeness

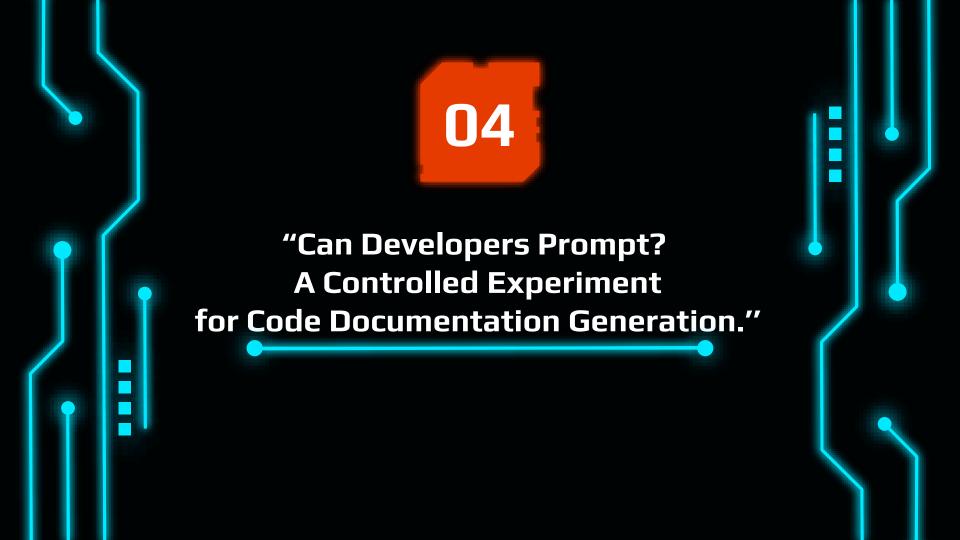
05 ROUGE-1

03 Conciseness

Results

- Document Level: GPT-4 is more accurate than the reference but slightly less complete.
- Class Level: GPT-4 performed poorly in completeness, correctness was high.
- Method-level: BLEU and ROUGE scores were significantly higher than other levels.

```
# checking response.status_come (if you on 5
 if response.status_code != 200:
       print(f"Status: {response.status_code - Dy moneus be
  else:
        print(f"Status: (response.status_code \ur)
 # using BeautifulSoup to parse the remove
 SOUP = BeautifulSoup(response.content, "nod.asser")
 images = soup.find_all("asg", attractates was seen
  # down toading image.
```





Idea

- Study how well both students and professionals, can prompt LLMs to generate code documentation
- compare ad-hoc prompts with predefined few-shot prompts.

- □ controlled, between-subject experiment
 - Used 2 Python functions to document
 - □ Ad-hoc group craft unique prompts
- Control group used pre-defined prompts

Evaluation metrics used:

01 Correctness

04 Usefulness

02 Missing Information

05 Helpfulness

03 Unnecessary Information

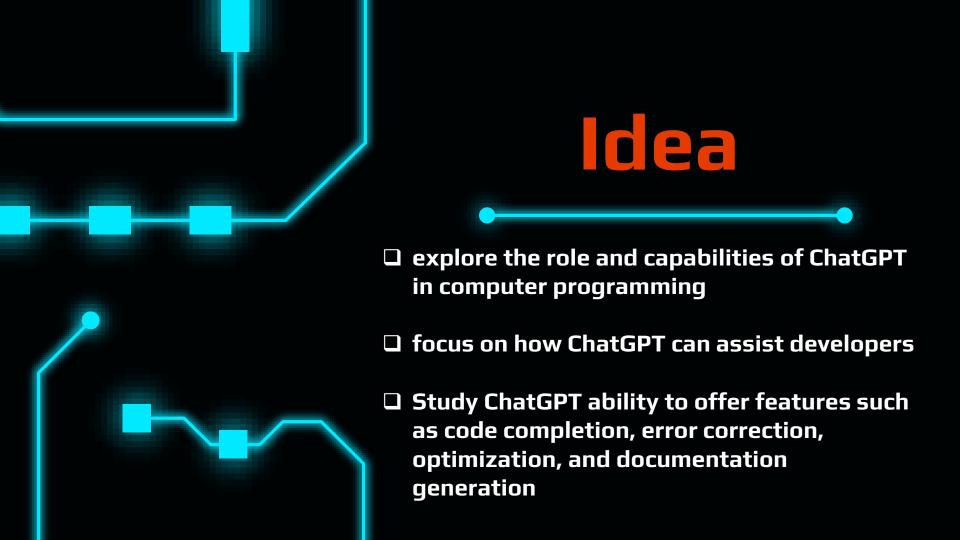
06 Readability

Results

- Predefined prompts produced higher-quality documentation.
- Professionals performed better with ad-hoc than students
- Both tools improved code understanding.
- Tools performed best with complex functions

```
# checking response.status_come (if you we
 if response.status_code != 200:
       print(f"Status: {response.status_code - by research
 else:
        print(f"Status: (response.status_code \ur)
 # using BeautifulSoup to parse DM resu
 SOUP = BeautifulSoup(response.content, "montenance")
images = soup.find_all("asg", attractates was seen
 # down toading image
```





- Giving ChatGPT prompts.
 - Analyzing its ability to:
 - 1) Predict code snippets
 - 2) Fix syntax errors
 - 3) Suggest optimization
- 4) Generate documentation

Evaluation Metrics

- No formal evaluation metrics are mentioned in the research.
- the effectiveness of ChatGPT is assessed qualitatively through examples
- Accuracy in code suggestions, error detection, and document generation.

Results

- □ ChatGPT is a valuable tool for developers, helps reduce coding errors, improve code quality, and accelerate development processes.
- ChatGPT is effective in providing automating documentation and refine code structure.
- ChatGPT still requires human supervision for tasks.

