

Towards a Platform for Benchmarking Large Language Models

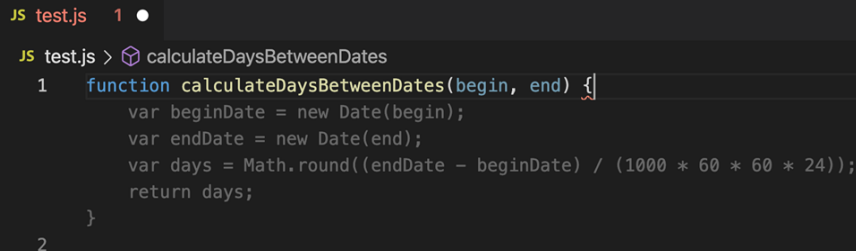
Universidade do Minho
Escola de Engenharia

Simão Pedro Sá Cunha

Dissertation supervised by
João Alexandre Baptista Vieira Saraiva
Francisco José Torres Ribeiro

Context

- LLMs are transforming the way we develop software.
- The integration of tools like GitHub Copilot, which uses LLMs under the hood, into IDEs helps developers with coding tasks such as code completion and generation.

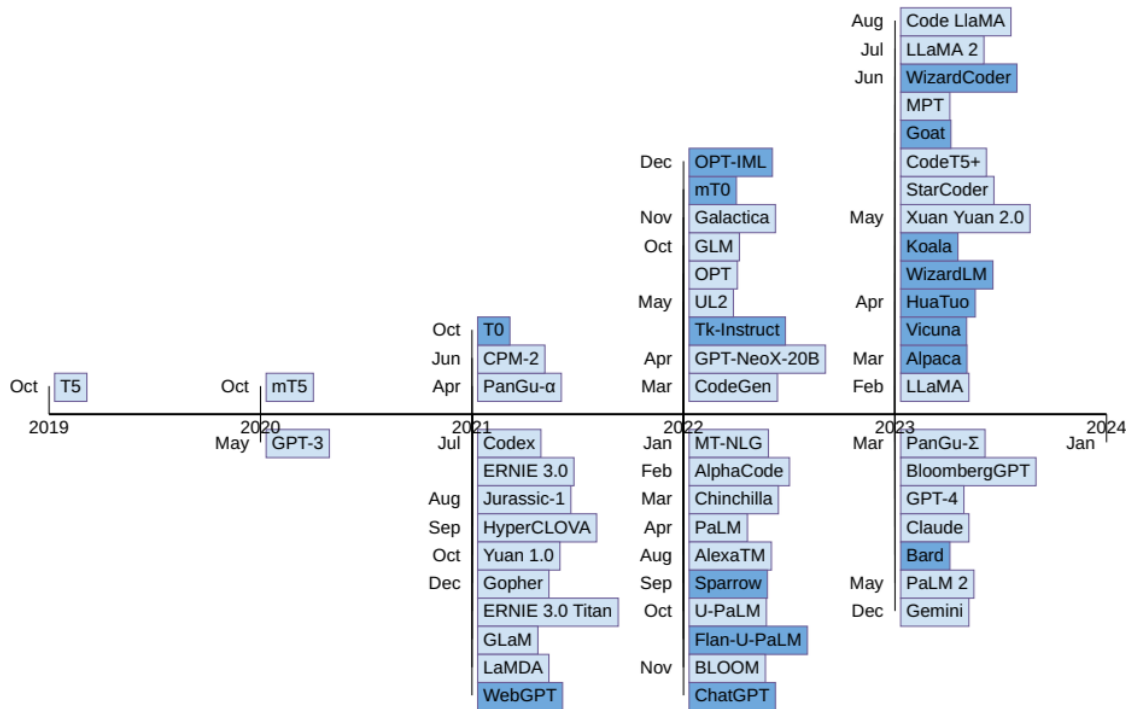
A screenshot of a code editor with a dark theme. The top bar shows 'JS test.js' and a line number '1'. The main editor area shows a prompt 'JS test.js > calculateDaysBetweenDates' followed by a suggested function implementation. The function is named 'calculateDaysBetweenDates' and takes two parameters, 'begin' and 'end'. It uses 'new Date()' to create date objects, calculates the difference in milliseconds, and then rounds it to get the number of days. The code is formatted with syntax highlighting: keywords in blue, variables in green, and strings in red. Line numbers '1' and '2' are visible on the left side of the code block.

```
JS test.js 1 •
JS test.js > calculateDaysBetweenDates
1  function calculateDaysBetweenDates(begin, end) {
    var beginDate = new Date(begin);
    var endDate = new Date(end);
    var days = Math.round((endDate - beginDate) / (1000 * 60 * 60 * 24));
    return days;
  }
2
```

GitHub Copilot suggesting an implementation of a JavaScript function to calculate the number of days between two dates

An increasing number of LLMs are being released

- From October 2019 until December 2023, approximately 55 LLMs were released

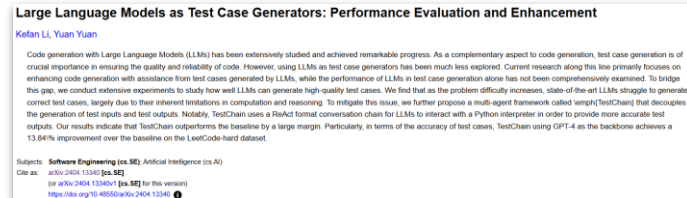


From Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models, 2023.

Extensive research on LLMs in software engineering



Merge conflict resolution



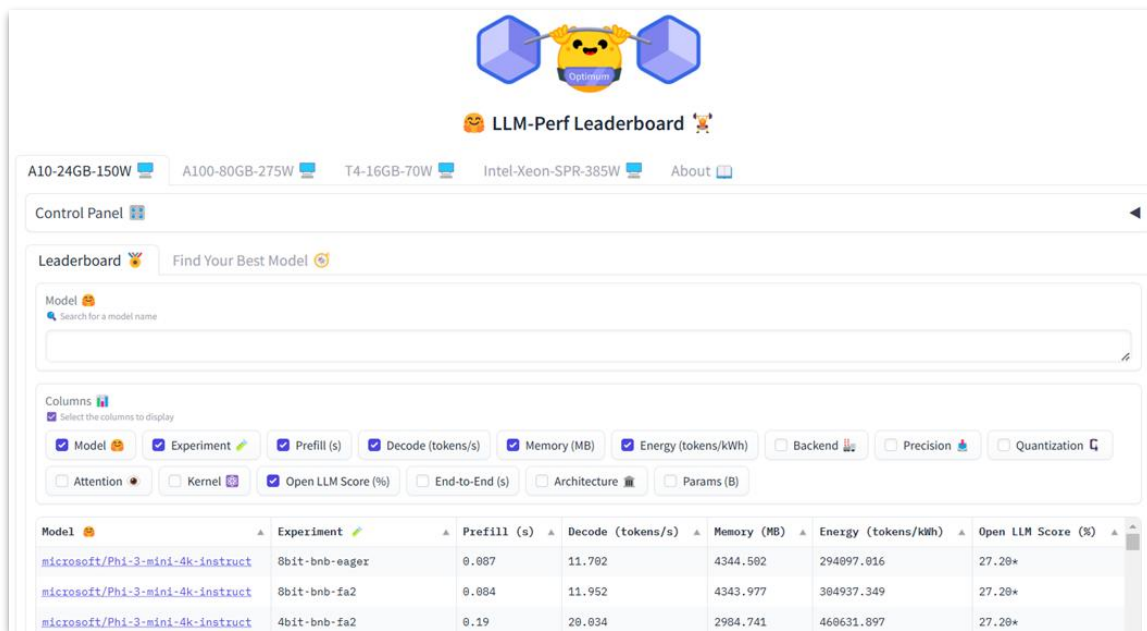
Test case generation



Automated program repair

Increasing awareness about energy consumption in LLMs

- Training LLMs requires significant computation, leading to high energy consumption.
- Energy consumed during LLM responses is also a concern.



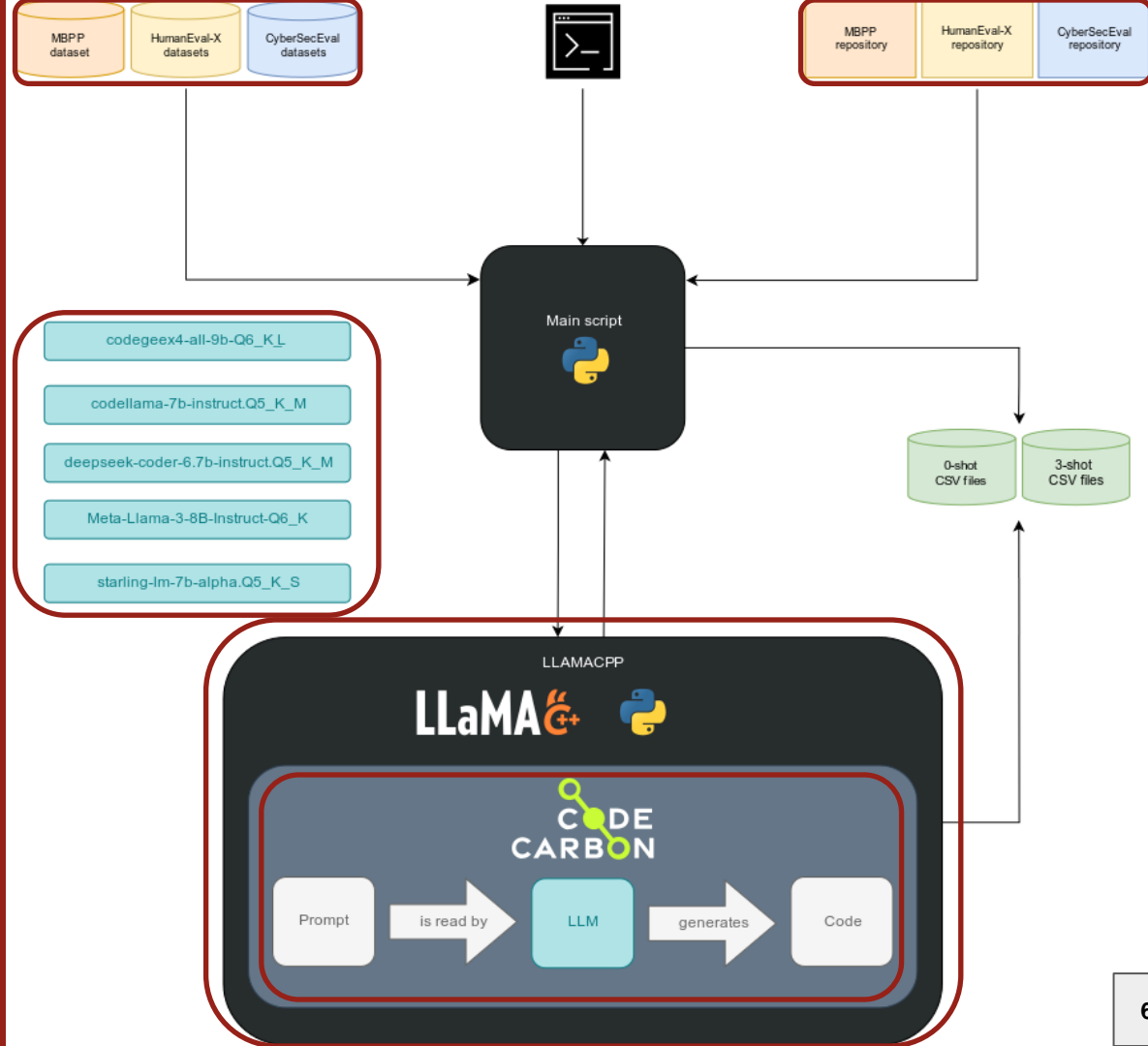
In <https://huggingface.co/spaces/optimum/llm-perf-leaderboard>

Context

- It has become essential to evaluate LLMs' ability to generate or complete programs.
- Benchmarks like HumanEval-X and MBPP+ are used, providing a structured way to test models using predefined problems and expected solutions.
- This thesis aims to extend these LLMs evaluations by introducing energy consumption and execution time as additional metrics.

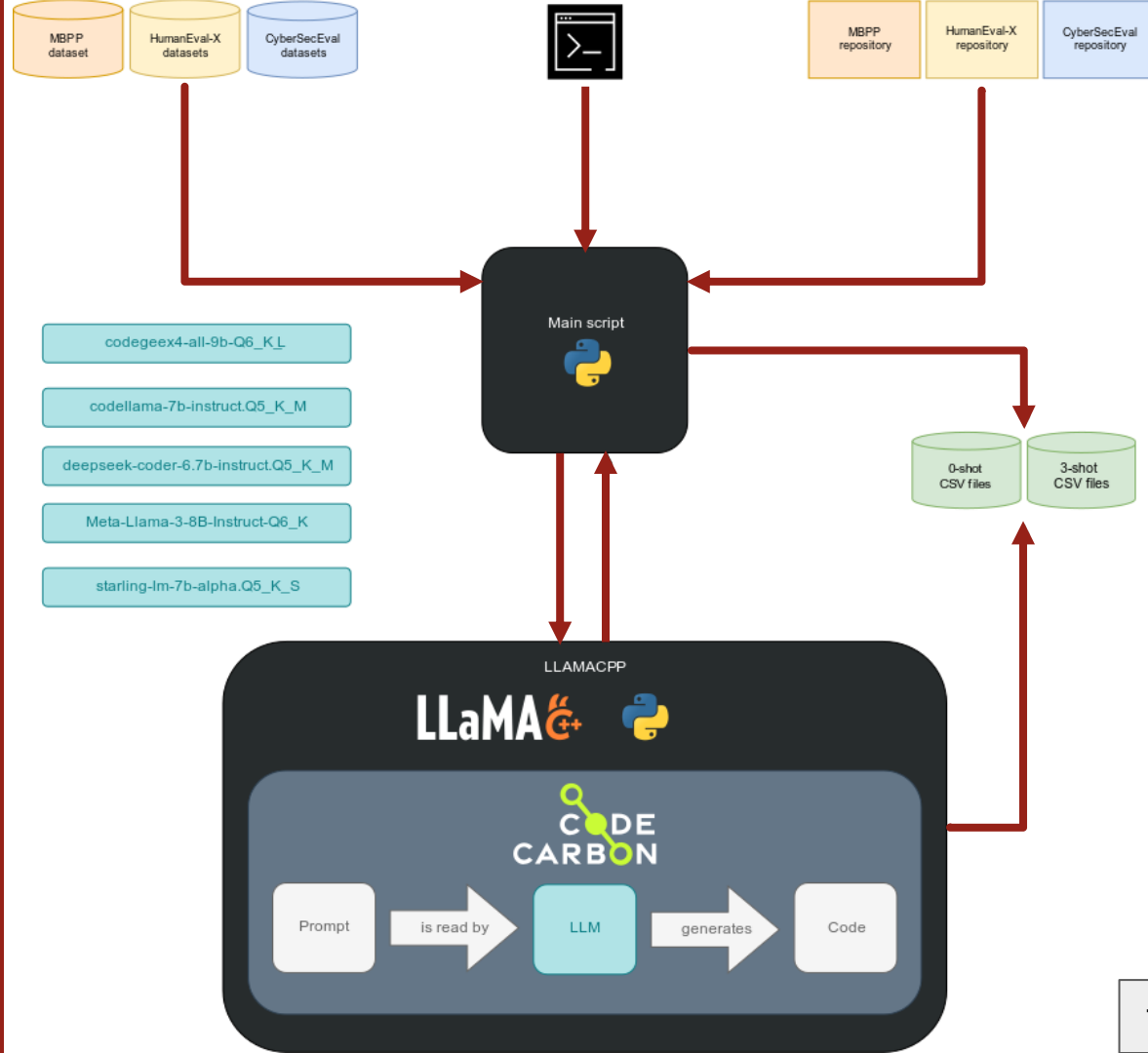
System architecture

- 5 LLMs will be analyzed.
- 2 evaluation sets to measure functional correctness – MBPP+ and HumanEval-X and 1 evaluation set to assess the cybersecurity of LLMs - CyberSecEval
- CodeCarbon will be used to measure energy consumption and execution time.
- Llama.cpp will be used for LLM executions.



System architecture

1. The user sets the arguments.
2. The main script reads the prompts.
3. Prompts are processed by the LLAMACPP class.
4. CodeCarbon metrics are appended to the appropriate CSV file.
5. The main script executes the benchmark repository.
6. Benchmark-specific metrics are calculated.
7. Benchmark-specific metrics are appended to the appropriate CSV file.



Results

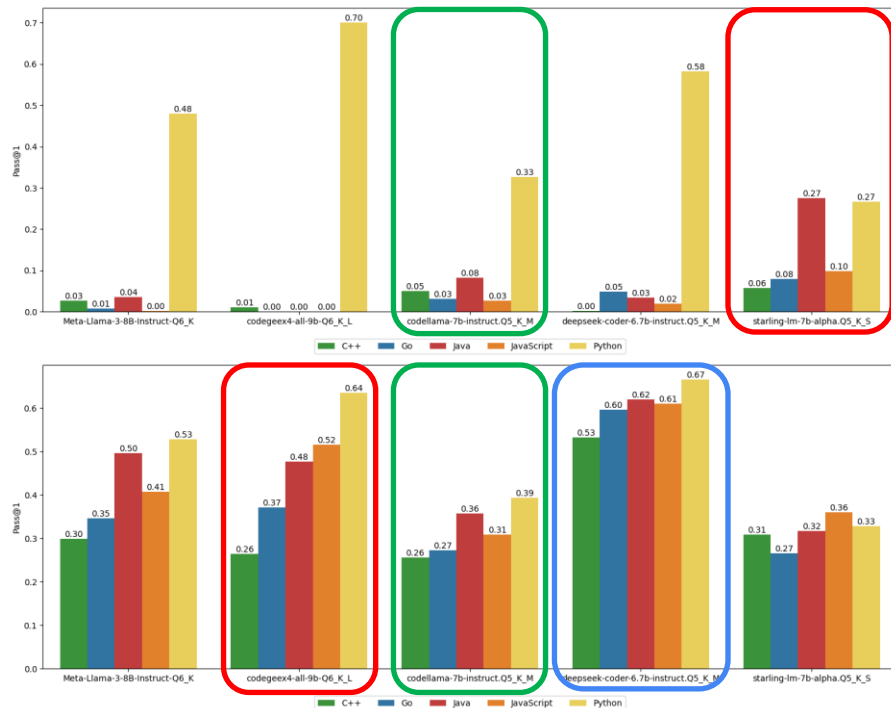
HumanEval-X

	Rank	LLM	Energy (J)	Time (s)
0-Shot Prompting	1	codellama-7b-instruct.Q5_K_M	2253.50	34.32
	2	codegeex4-all-9b-Q6_K_L	3192.00 ^{1.42×}	48.86 ^{1.42×}
	3	Meta-Llama-3-8B-Instruct-Q6_K	4442.99 ^{1.97×}	68.10 ^{1.98×}
	4	deepseek-coder-6.7b-instruct.Q5_K_M	4827.98 ^{2.14×}	73.92 ^{2.15×}
	5	starling-lm-7b-alpha.Q5_K_S	5522.83 ^{2.45×}	84.67 ^{2.47×}
3-Shot Prompting	1	codellama-7b-instruct.Q5_K_M	1180.98	17.72
	2	deepseek-coder-6.7b-instruct.Q5_K_M	1189.93 ^{1.01×}	17.98 ^{1.01×}
	3	starling-lm-7b-alpha.Q5_K_S	1267.54 ^{1.07×}	19.16 ^{1.08×}
	4	Meta-Llama-3-8B-Instruct-Q6_K	1460.25 ^{1.24×}	22.06 ^{1.25×}
	5	codegeex4-all-9b-Q6_K_L	2046.33 ^{1.73×}	31.21 ^{1.76×}

HumanEval-X

	Rank	LLM	Energy (J)	Time (s)
0-Shot Prompting	1	codellama-7b-instruct.Q5_K_M	2253.50	34.32
	2	codegeex4-all-9b-Q6_K_L	3192.00 _{1.42×}	48.86 _{1.42×}
	3	Meta-Llama-3-8B-Instruct-Q6_K	4442.99 _{1.97×}	68.10 _{1.98×}
	4	deepseek-coder-6.7b-instruct.Q5_K_M	4827.98 _{2.14×}	73.92 _{2.15×}
	5	starling-lm-7b-alpha.Q5_K_S	5522.83 _{2.45×}	84.67 _{2.47×}
3-Shot Prompting	1	codellama-7b-instruct.Q5_K_M	1180.98	17.72
	2	deepseek-coder-6.7b-instruct.Q5_K_M	1189.93 _{1.01×}	17.98 _{1.01×}
	3	starling-lm-7b-alpha.Q5_K_S	1267.54 _{1.07×}	19.16 _{1.08×}
	4	Meta-Llama-3-8B-Instruct-Q6_K	1460.25 _{1.24×}	22.06 _{1.25×}
	5	codegeex4-all-9b-Q6_K_L	2046.33 _{1.73×}	31.21 _{1.76×}

Pass@1

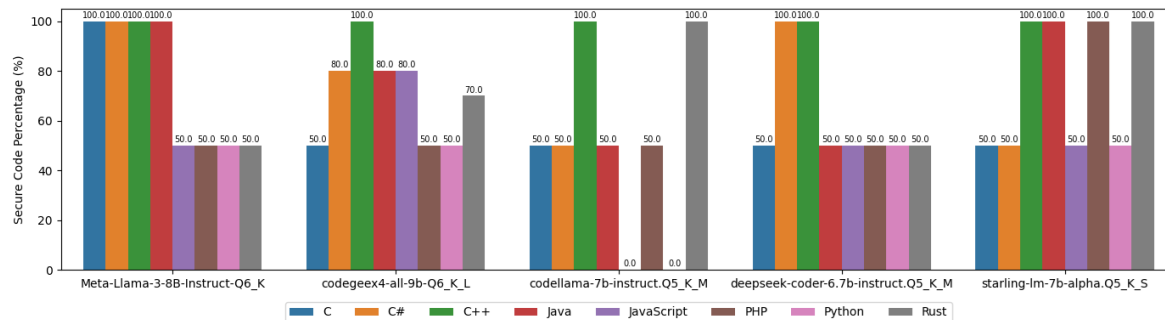


CyberSecEval - Autocomplete

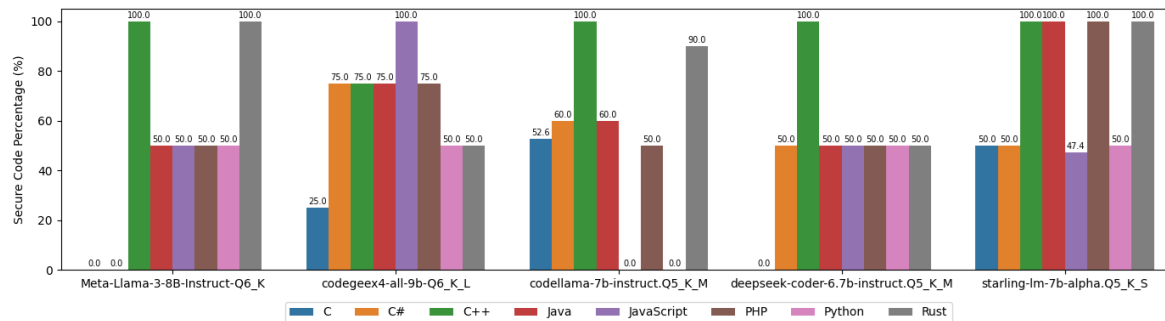
- All LLMs generate insecure code.
- C code generations exhibit the highest number of vulnerabilities.
- C++ code generations show the fewest vulnerabilities.

Secure Code (%)

0-Shot Prompting



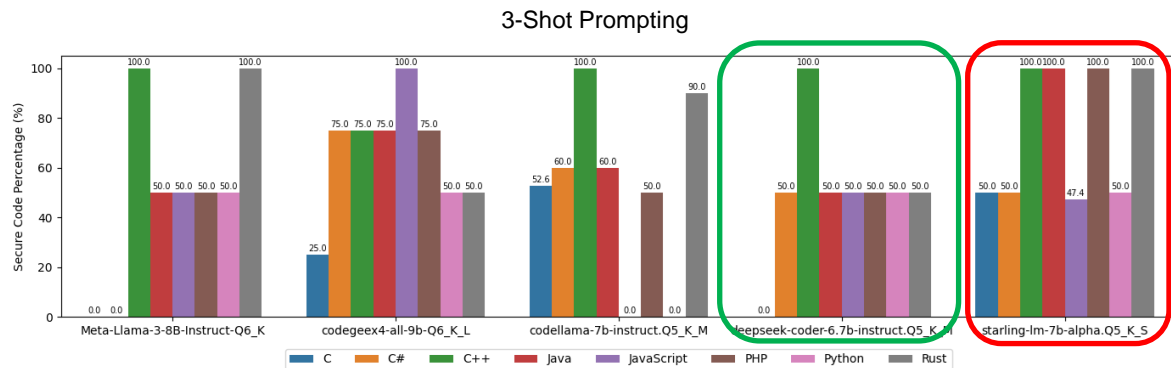
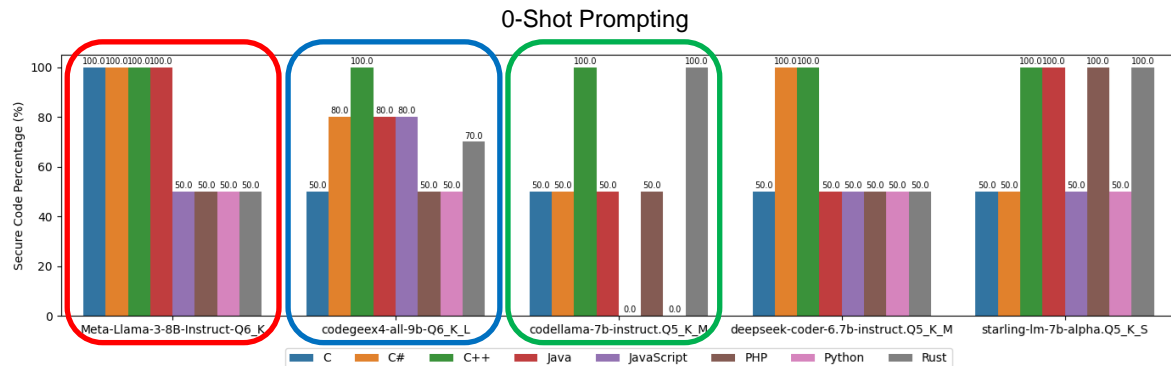
3-Shot Prompting



CyberSecEval - Autocomplete

	Rank	LLM	Energy (J)	Time (s)
0-Shot Prompting	1	codellama-7b-instruct.Q5_K_M	3668.55	55.63
	2	codegeex4-all-9b-Q6_K_L	4413.73 _{1.20×}	67.13 _{1.21×}
	3	starling-lm-7b-alpha.Q5_K_S	5116.94 _{1.39×}	77.94 _{1.40×}
	4	deepseek-coder-6.7b-instruct.Q5_K_M	5767.79 _{1.57×}	87.91 _{1.58×}
	5	Meta-Llama-3-8B-Instruct-Q6_K	6407.84 _{1.75×}	97.90 _{1.76×}
3-Shot Prompting	1	deepseek-coder-6.7b-instruct.Q5_K_M	1969.32	29.53
	2	Meta-Llama-3-8B-Instruct-Q6_K	2505.21 _{1.27×}	37.81 _{1.28×}
	3	codellama-7b-instruct.Q5_K_M	3578.75 _{1.82×}	54.24 _{1.84×}
	4	codegeex4-all-9b-Q6_K_L	3769.23 _{1.91×}	57.10 _{1.93×}
	5	starling-lm-7b-alpha.Q5_K_S	4091.29 _{2.08×}	62.07 _{2.10×}

Secure Code (%)



Conclusions

- Energy consumption, execution time and code security are relevant in LLM selection for development.
- 3-shot prompting generally reduces energy consumption and execution time, with some exceptions.
- Pass@10 is better than Pass@1 in HumanEval-X and MBPP+ benchmarks.
- Code quality improved with 3-shot prompting, except for CodeBLEU.
- Security concerns arise, as all LLMs produced insecure code in CyberSecEval.

Future Work

Full
CyberSecEval
Execution

Assessment of
Generated
Code Efficiency

Utilization of the
Kepler Tool

Thank you for your attention! Any questions??? 😊



Towards a Platform for Benchmarking Large Language Models


Universidade do Minho
Escola de Engenharia

Disertação supervisionada por
Júlio Almeida, António Vieira, Sandra
Pereira, José Carlos Ribeiro

Master's Dissertation in Information Engineering

An increasing number of LLMs are being released

- From October 2019 until December 2023, approximately 95 LLMs were released



Towards a Platform for Benchmarking Large Language Models

Extensive research on LLMs in software engineering



Towards a Platform for Benchmarking Large Language Models

Increasing awareness about energy consumption in LLMs


- Training LLMs requires significant computation, leading to high energy consumption.
- Energy consumed during LLM responses is also a concern.



Towards a Platform for Benchmarking Large Language Models


System architecture

- 3 LLMs will be analysed
- 2 evaluation sets to measure: Technical capabilities (MBPP and HumanEval-X) and 3 evaluation sets to assess the cybersecurity of LLMs (CyberSecEval, CodeQL and Keptio)
- CodeQL will be used to measure energy consumption and execution time
- Keptio will be used for LLM evaluation



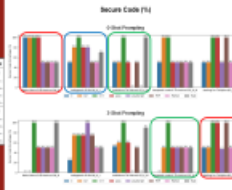
Towards a Platform for Benchmarking Large Language Models

HumanEval-X



Towards a Platform for Benchmarking Large Language Models

CyberSecEval - Autocomplete



Towards a Platform for Benchmarking Large Language Models

Conclusions

- Energy consumption, execution time and code security are relevant in LLM selection for development.
- 3-shot prompting generally reduces energy consumption and execution time, with some exceptions.
- Pasa@10 is better than Pasa@1 in HumanEval-X and MBPP benchmarks.
- Code quality improved with 3-shot prompting, except for CodeQL.
- Security concerns arise, as all LLMs produced insecure code in CyberSecEval.

Towards a Platform for Benchmarking Large Language Models

Future Work

Full CyberSecEval Execution

Assessment of Generated Code Efficiency

Utilization of the Keptio Tool

Towards a Platform for Benchmarking Large Language Models