# Code Llama: Open Foundation Models for Code

# 28m+ developers in the world



Increase productivity — 32.81%
Speed up learning — 25.17%
Greater efficiency — 24.96%
Improve accuracy in coding — 13.31%
Improve collaboration — 3.75%

# Code generation

- Program synthesis is not new, *including DL based*, e.g. (DeepCoder, Balog et al., 2016), (Bošnjak et al., 2017).
- Renewed interests with LLMs (e.g. Codex, 2021)

Tasks:

- Code completion
- Program synthesis from input/output pairs
- Linting
- Typing
- Bug finding
- Tests generation
- Translation

# Existing models

Closed models
- Running on GPUs on servers
- Inaccessible model weights



Open models (Llama, StarCoder)
- Can be finetuned for particular language/ codebases
- Can run locally, with no internet connection
- Benefit from community improvements
- The models are **free** to use, no license fee

# LLM 101

- P(next token | all previous tokens)
  - For all tokens
  - For **lots of text**
- `<bos> the cat sat on  the mat <eos>`
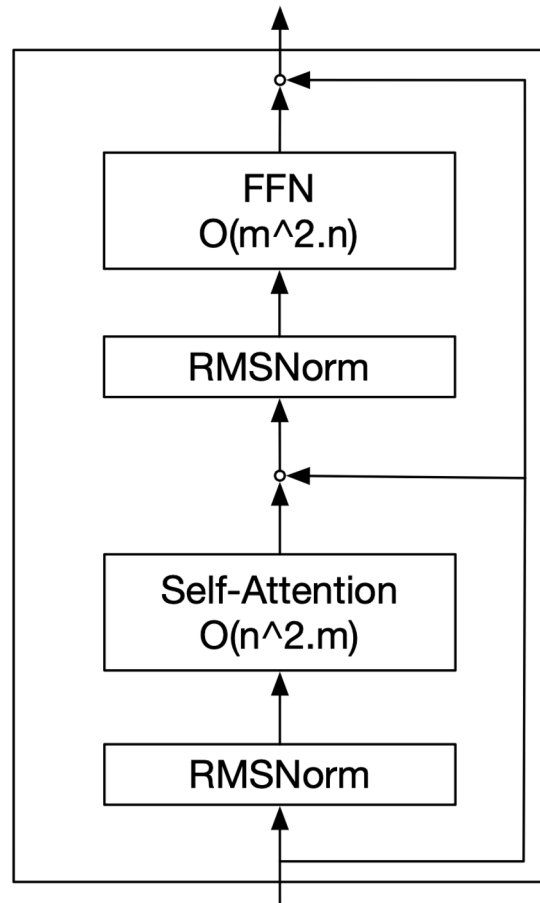  Teacher forcing:

```
<bos> ???
<bos> the ???
<bos> the cat ???
<bos> the cat sat ???
<bos> the cat sat on  ???
<bos> the cat sat on  the ???
<bos> the cat sat on  the mat ???
```

FFN
$O(m^2.n)$

RMSNorm

Self-Attention
$O(n^2.m)$

RMSNorm

# Code Llama



**PROMPT**

**RESPONSE**

Clear    Submit

∞ Meta AI

# Generating Code Llama's paper figures with Code Llama



Figure 3: **Correlations between Languages.** Correlation scores between the Python, C++, Java, PHP, C#, TypeScript (TS), and Bash, reported for different model sizes. The code for this figure was generated by CODE LLAMA - INSTRUCT, the prompt and code can be seen in Figure 22.

Llama 2

Code Llama

# Code Llama

500B tokens

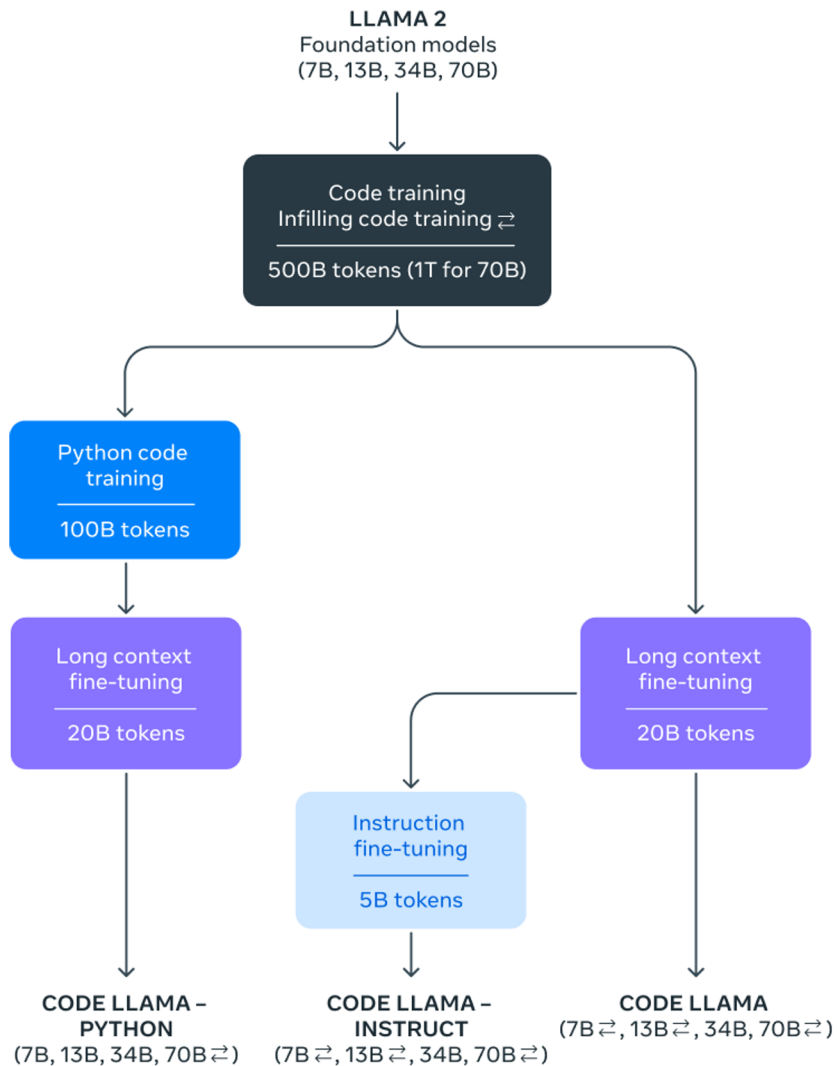Mostly programming languages
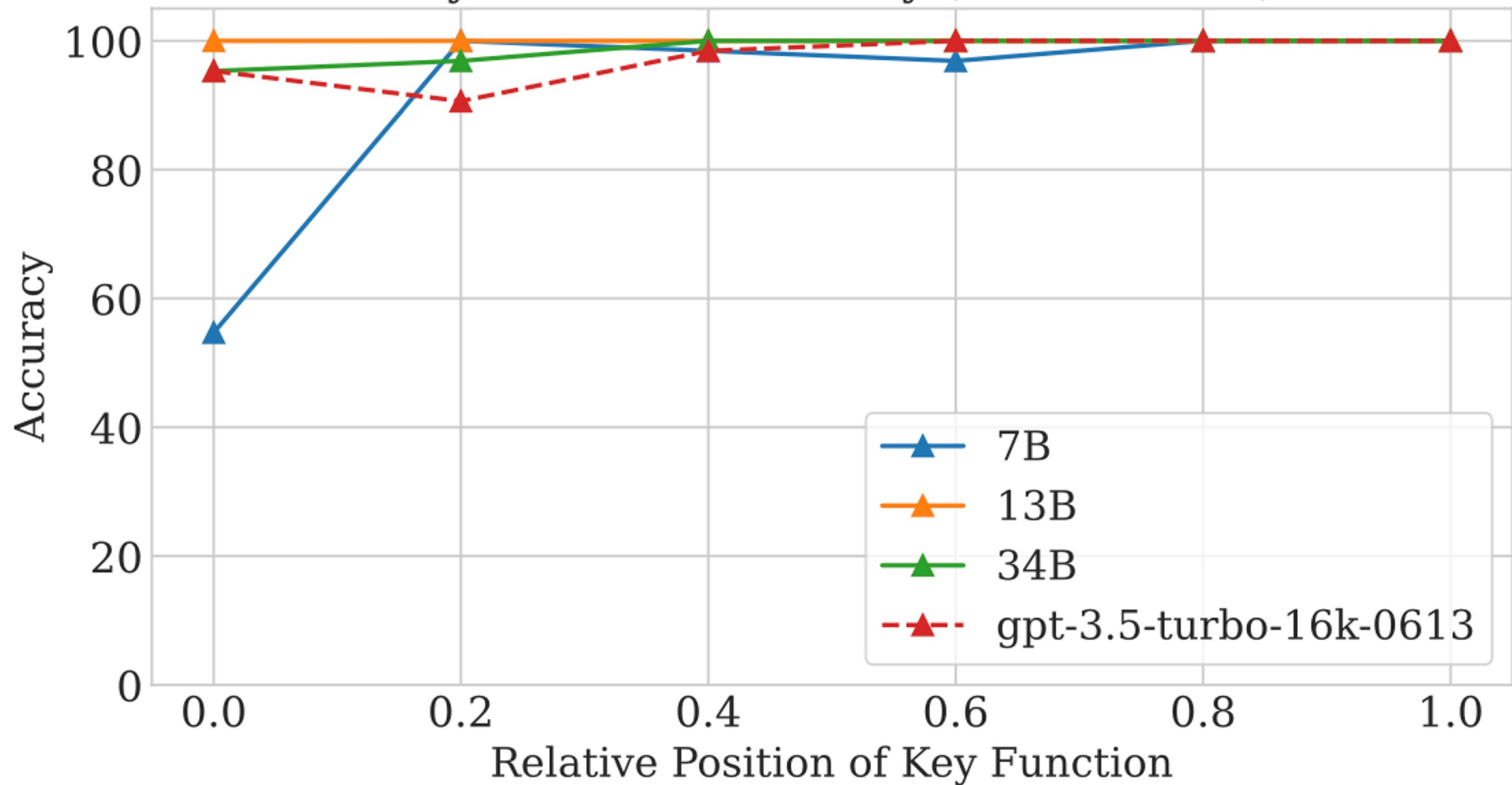
7B

13B

34B

70B

7B 13B 34B 70B

7B 13B 34B 70B

# Long context

- ~20B tokens fine-tuning
- Trained with up 16k tokens
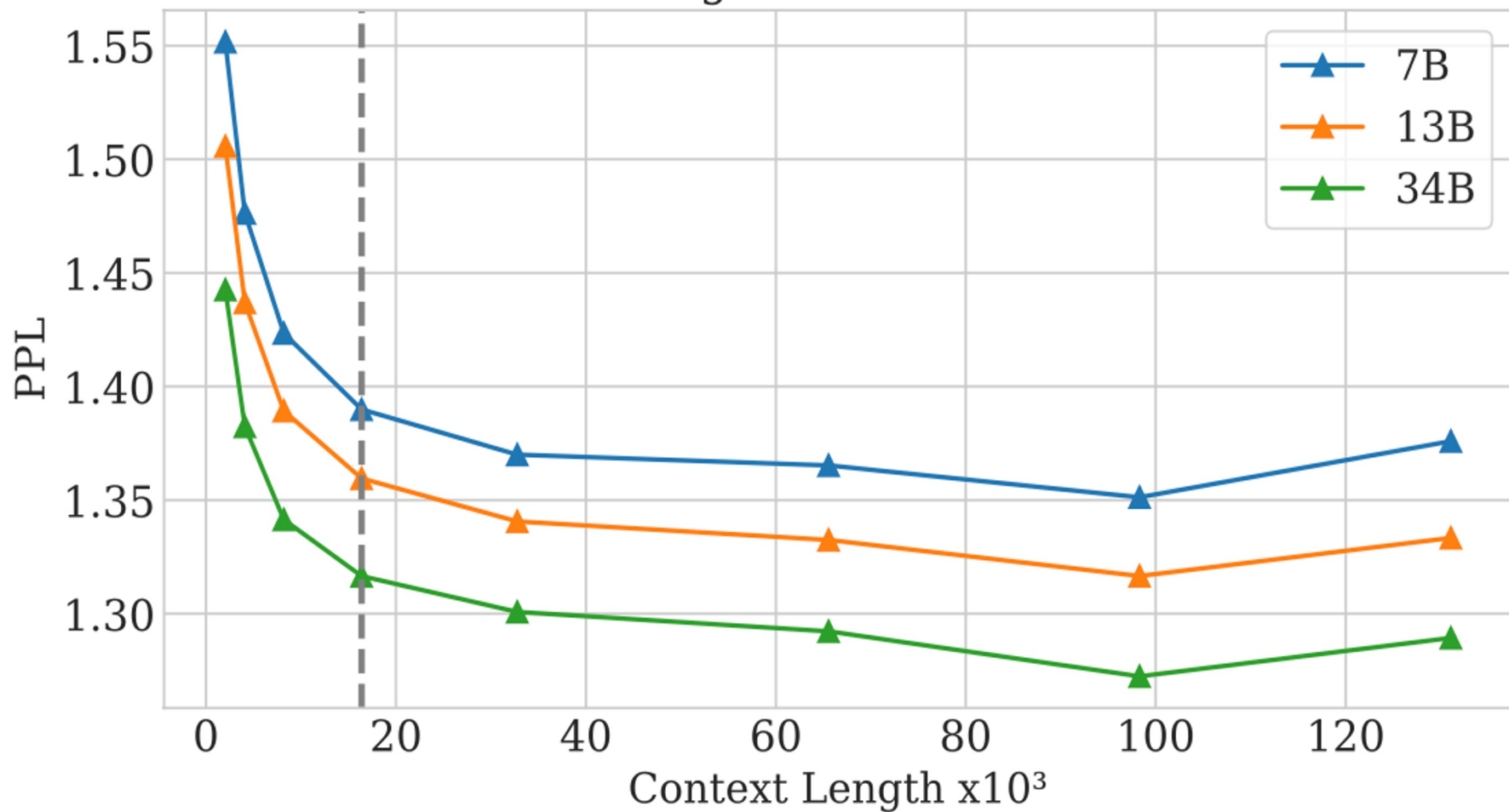- Supports up to 100k tokens = 8k lines of code
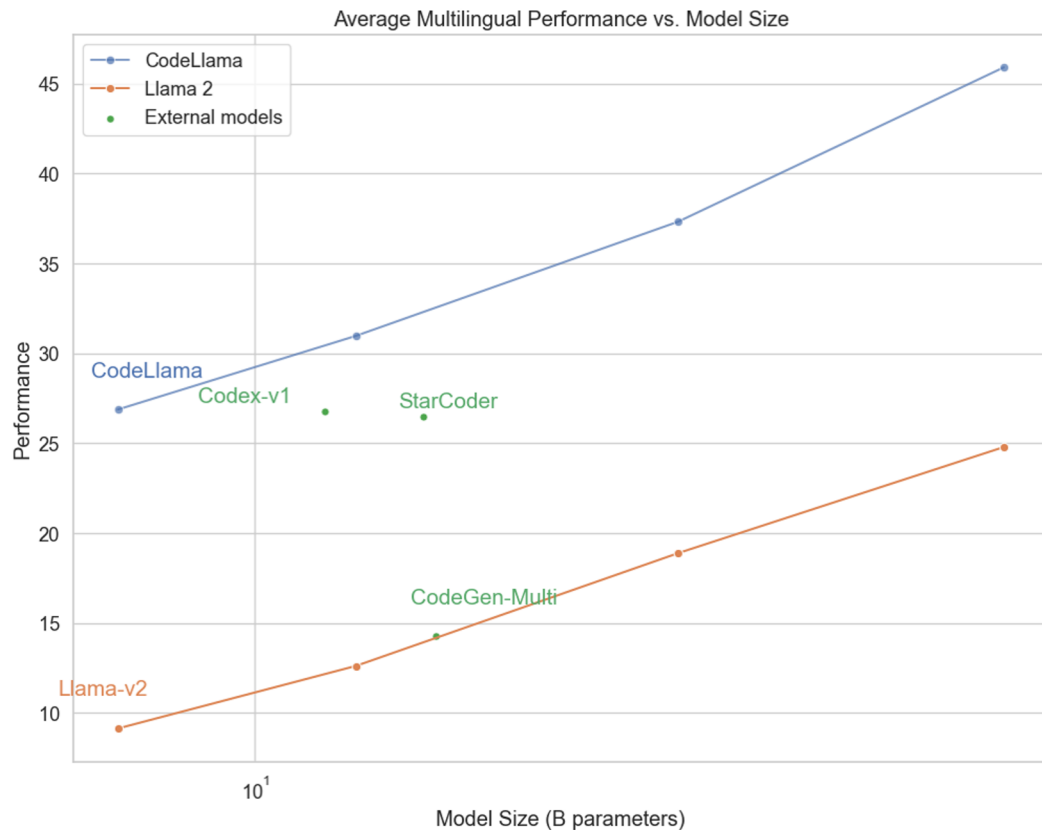
Key Retrieval Accuracy (~16K tokens)

Large Source Files

# Fill-in-the-middle (FIM)

```python
class Character:.py 1

Users > broz > workspaces > CodeLlama_autocomplete_tests > class Character:.py > ...
1  if __name__ == "__main__":
2      alice = Character("Alice", 45, "Software Engineer")
```
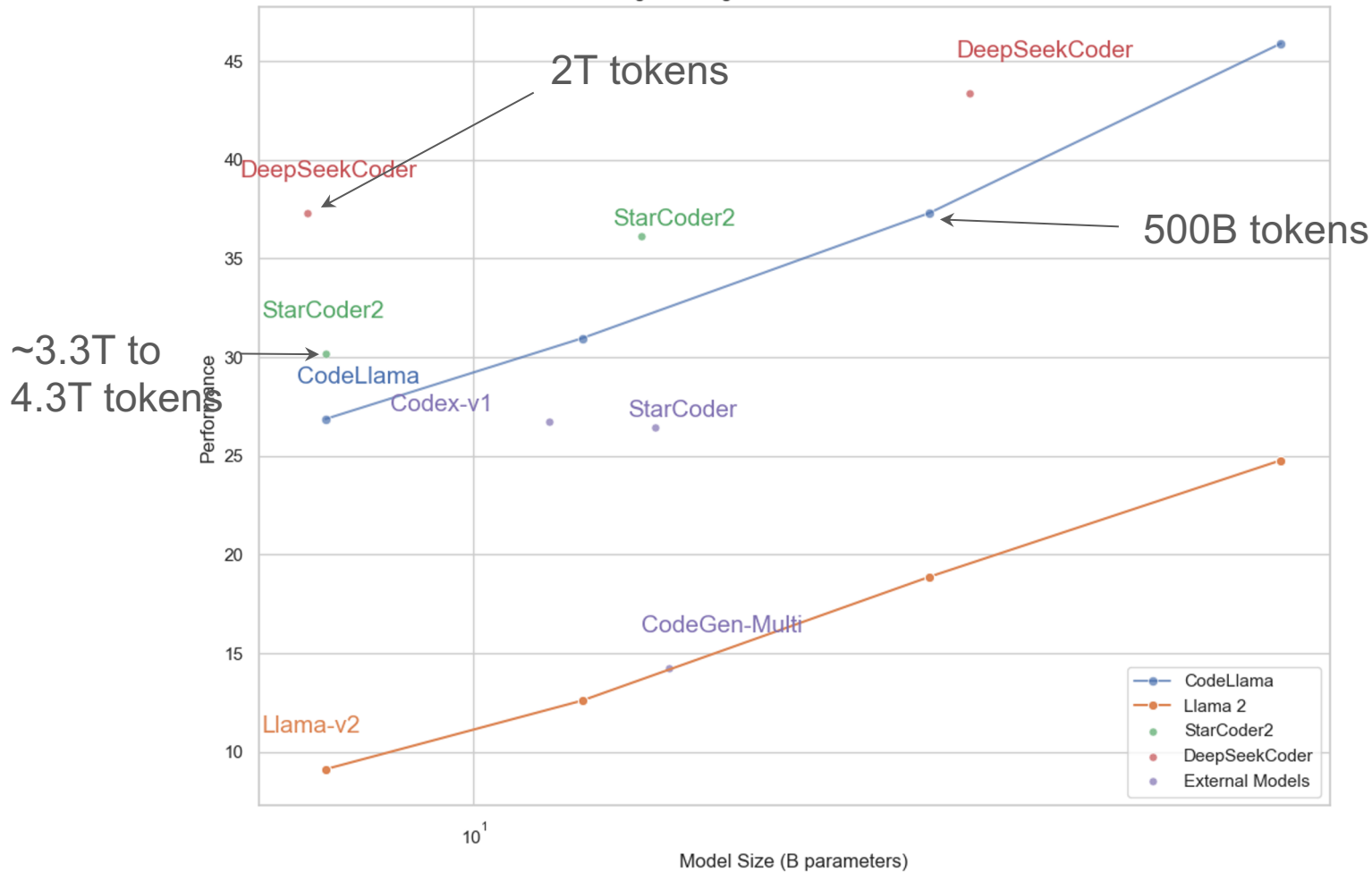
# HumanEval example

```python
def unique(l: list):
    """Return sorted unique elements in a list
    >>> unique([5, 3, 5, 2, 3, 3, 9, 0, 123])
    [0, 2, 3, 5, 9, 123]
    """
    return sorted(list(set(l)))
```
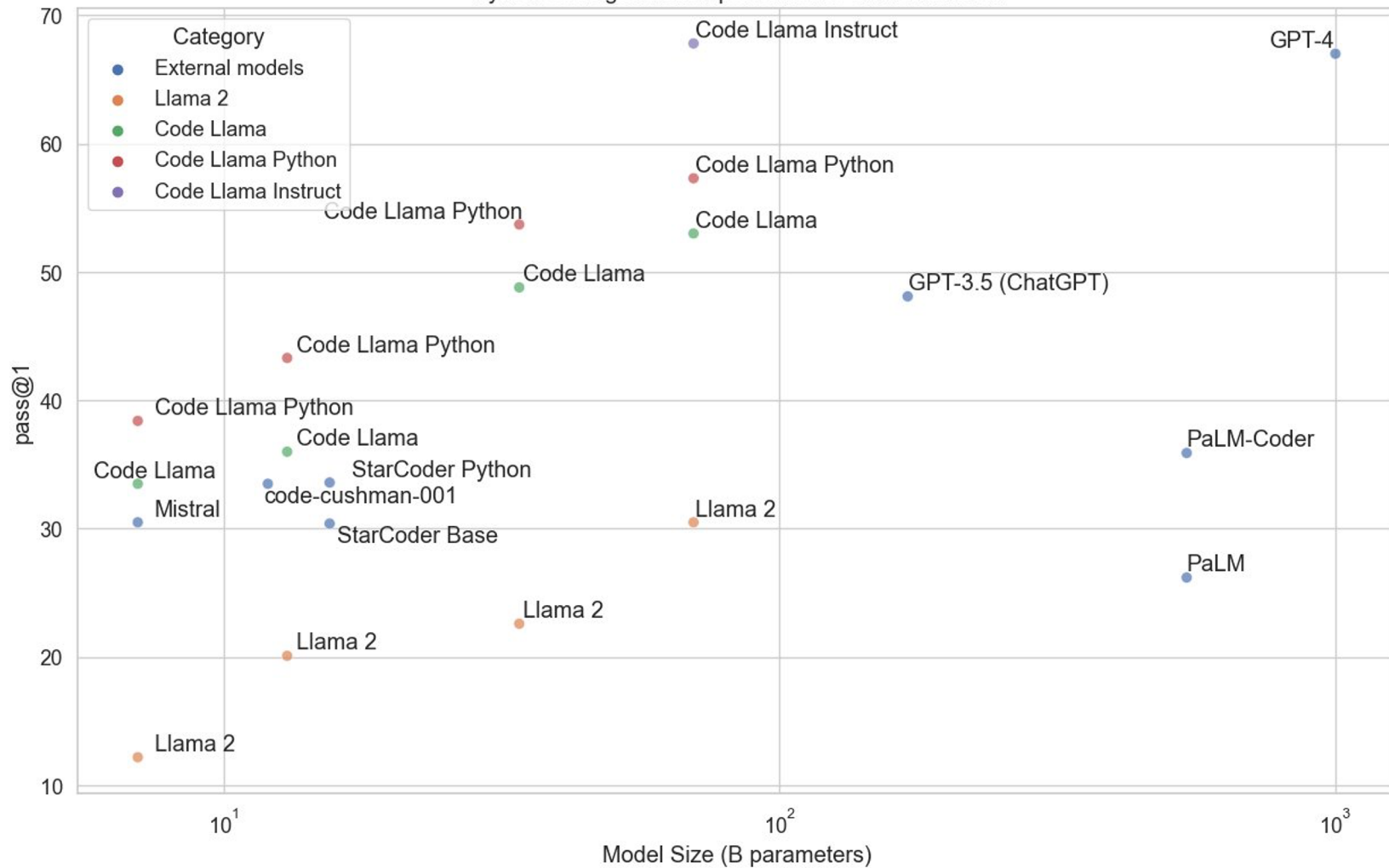
# Code Llama 7B outperforms Llama 2 70B on multilingual coding benchmarks



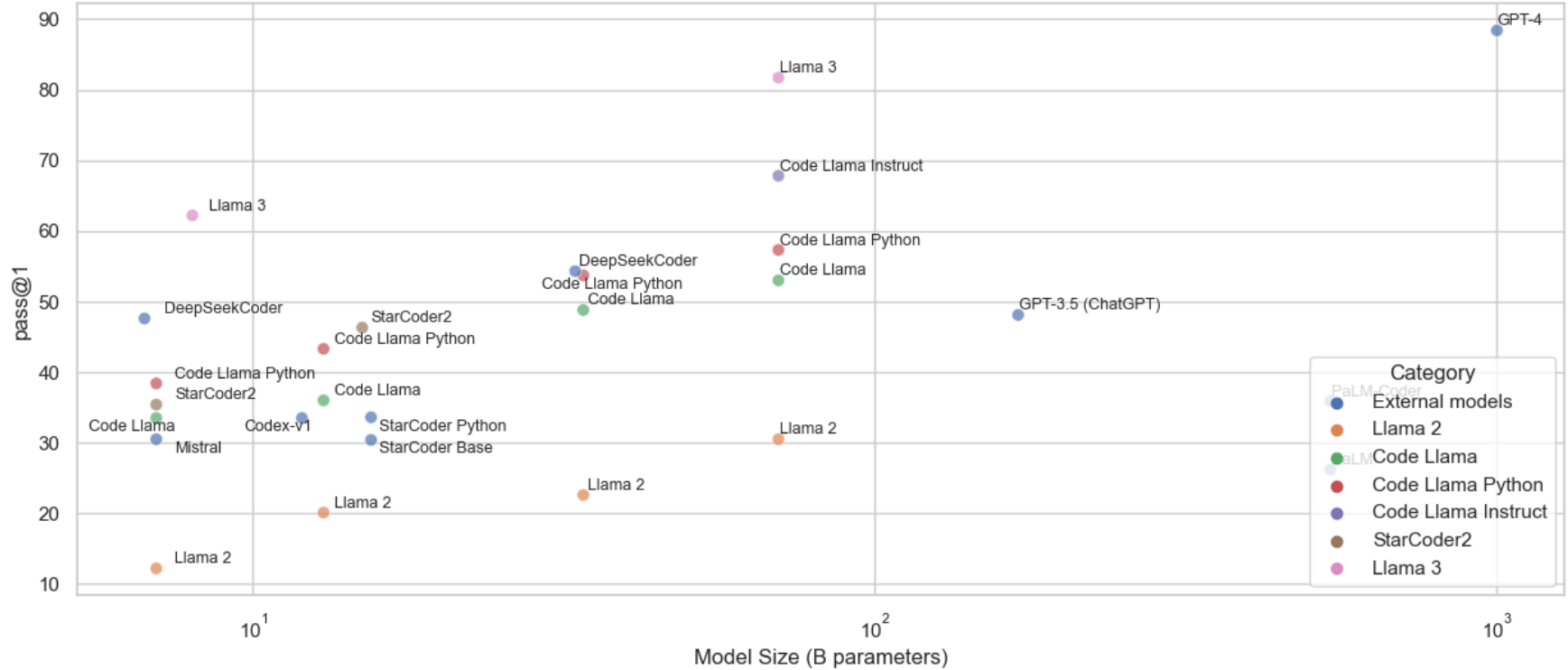Average Multilingual Performance vs. Model Size
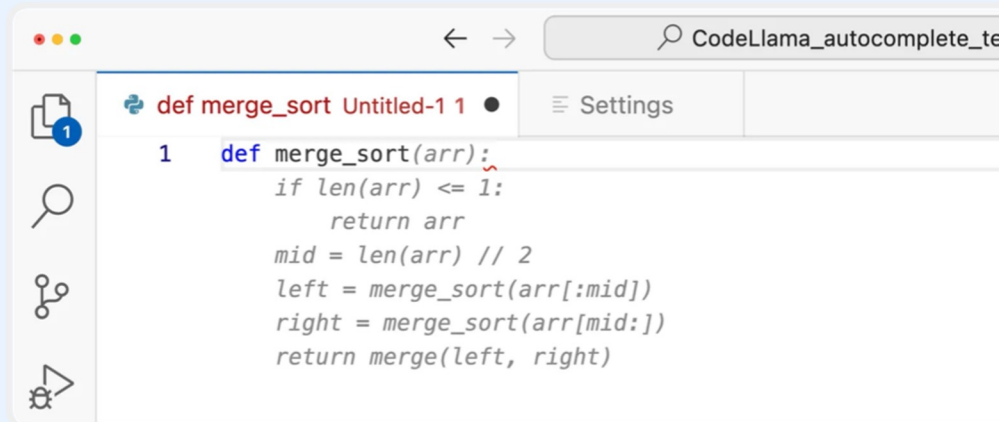
Average Multilingual Performance vs. Model Size

Python code generation performance on HumanEval

HumanEval Python Performance vs Model Size

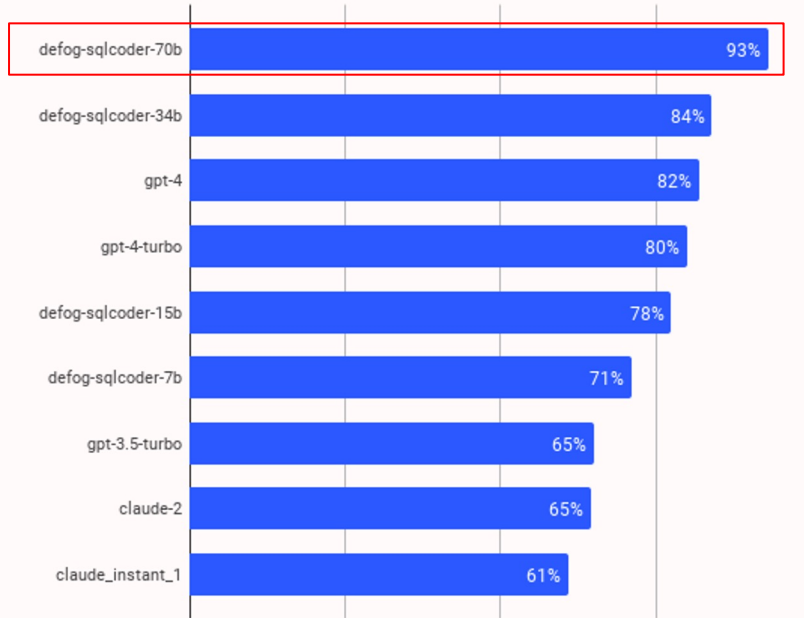# Hugging Face integration with VSCode

# Built on Code Llama, 2 examples:



Percentage of correctly generated SQL queries on novel schemas not seen in training (n = 200) in SQL-Eval

| | |
|---|---|
| defog-sqlcoder-70b | 93% |
| defog-sqlcoder-34b | 84% |
| gpt-4 | 82% |
| gpt-4-turbo | 80% |
| defog-sqlcoder-15b | 78% |
| defog-sqlcoder-7b | 71% |
| gpt-3.5-turbo | 65% |
| claude-2 | 65% |
| claude_instant_1 | 61% |

HumanEval

https://twitter.com/rishdotblog/status/1752329471867371659 20k instructions || https://www.phind.com/blog/introducing-phind-70b 50B tokens     27

Write a simple version of pong using pygame.

# Get started with Code Llama

- Ollama https://ollama.com/library/codellama
- HuggingFace https://huggingface.co/codellama/
- Perplexity AI chat https://labs.perplexity.ai/
- Our inference GitHub repository
  https://github.com/facebookresearch/codellama

Questions ?