

LLM-Supported Natural Language to Bash Translation

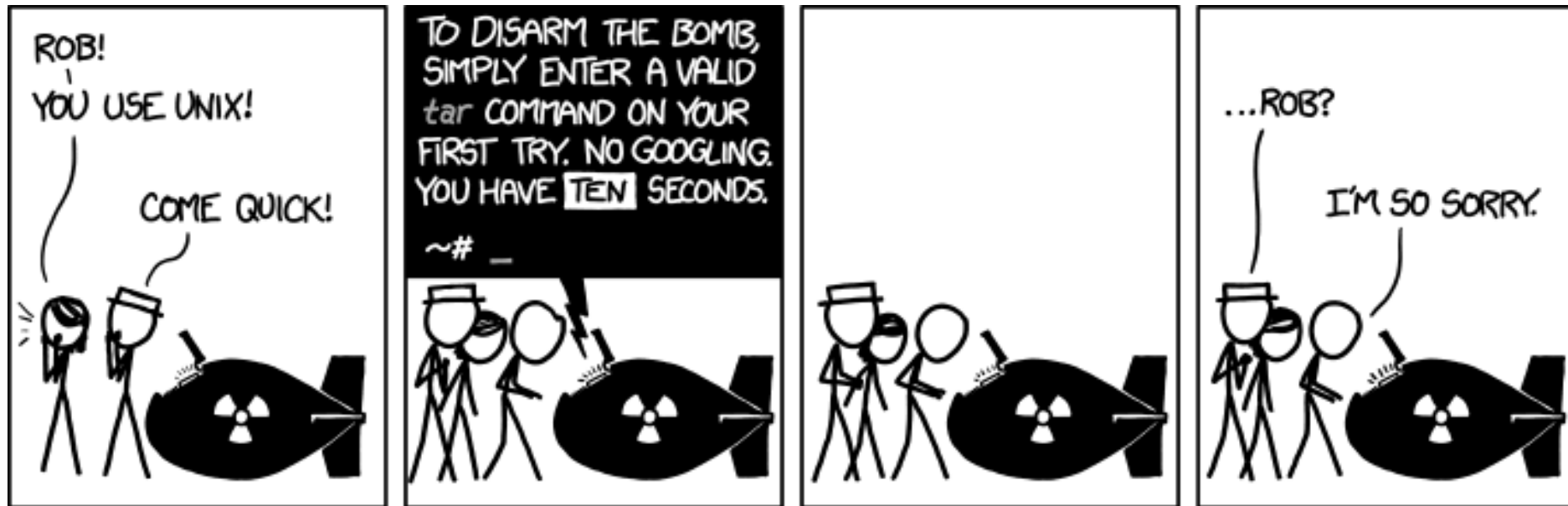
Finnian Westenfelder^{1,2}, Erik Hemberg¹, Miguel Tulla¹,
Stephen Moskal¹, Una-May O'Reilly¹, Silviu Chiricescu³

¹ALFA Group MIT-CSAIL, ²Draper Scholar, ³Charles Stark Draper Laboratory



Motivation

- Bash is the default Linux command line interface
- Bash commands are complex and difficult to memorize
- Incorrect commands can cause system failures



<https://xkcd.com/1168/>

Motivation

- The translation capabilities of Large Language Models (LLMs) can simplify command line interfaces
- Current LLMs are unreliable for translating natural language to Bash commands (NL2SH)
- Improving NL2SH requires
 - (1) clean data
 - (2) reliable benchmarks
 - (3) accurate translation methods

Input: Natural Language

List files in the /workspace directory that were accessed over an hour ago.

Output: Bash Command

`find /workspace -type f -amin +60`

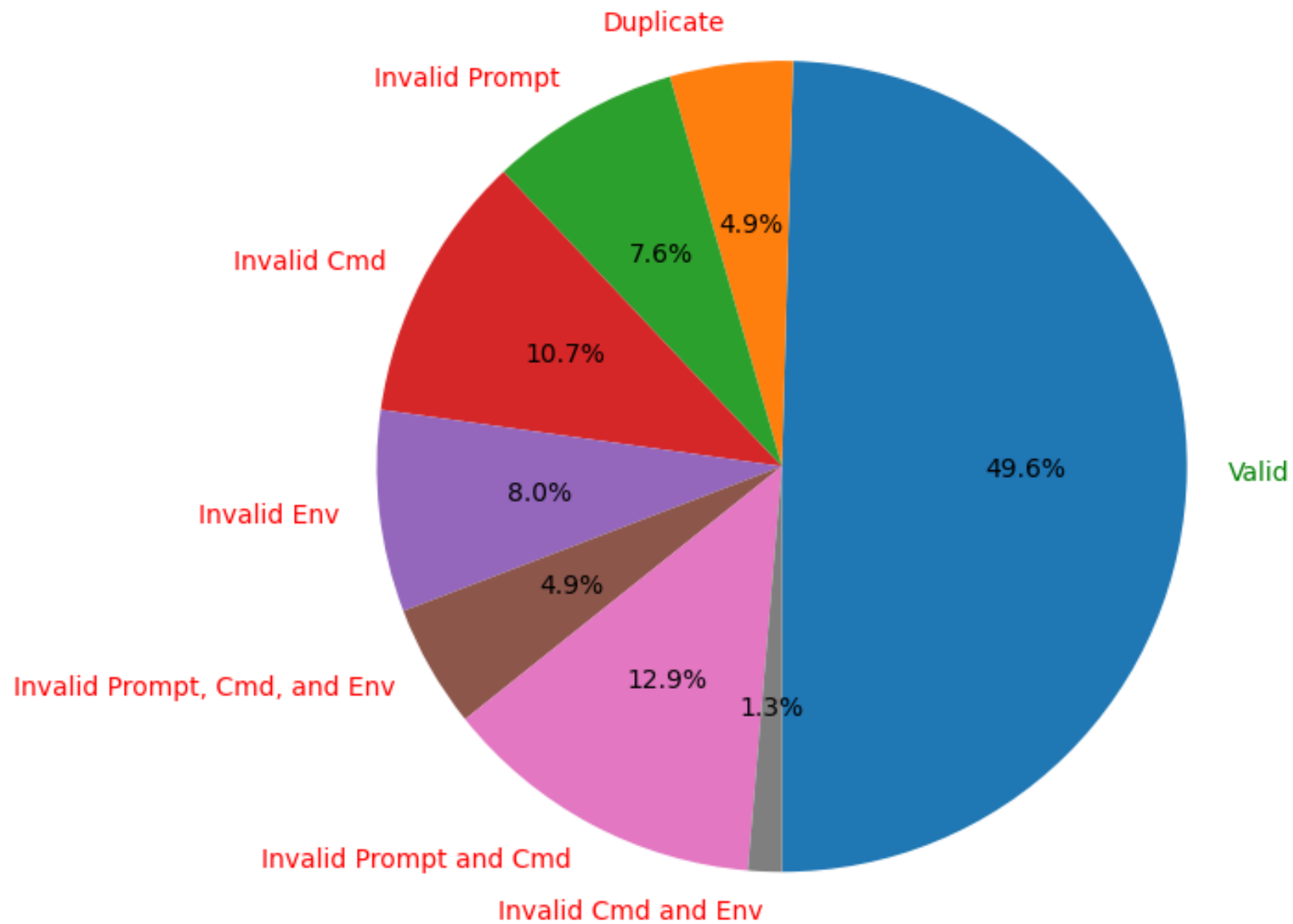
Dataset Overview

- **Goal:**

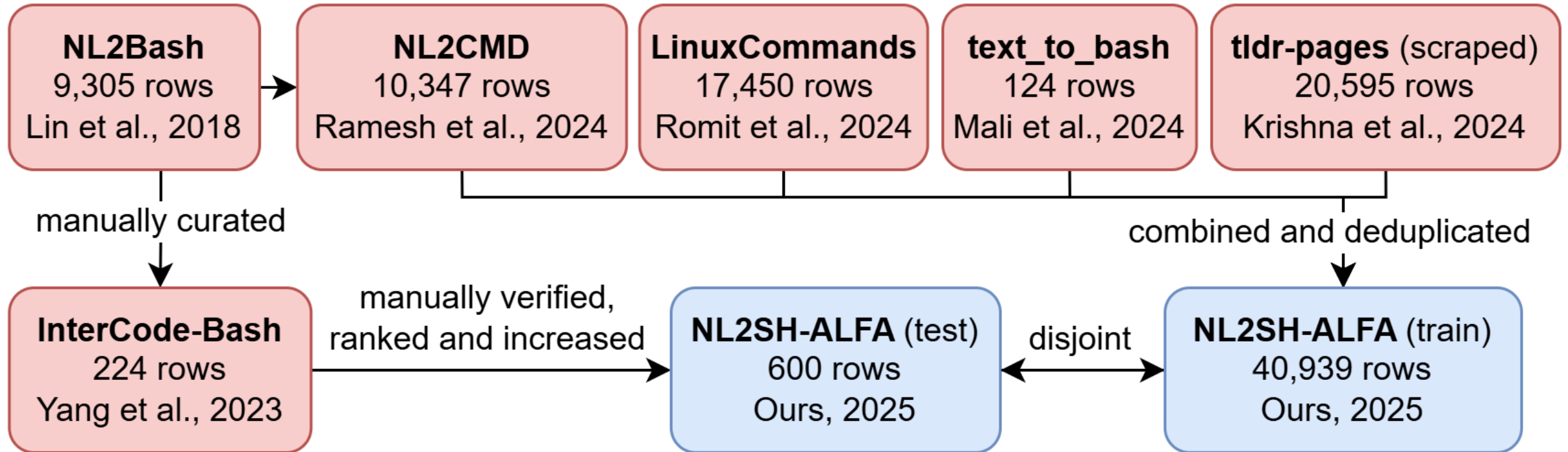
Given		Determine
Natural Language	Bash	Correct
print the system disk usage	df -h	True
remove a directory named foo	rm foo	False
print the current user's id	id -u	True

- **Challenge:** Datasets contain invalid translations.
- **Research Question:** How can we validate NL2SH datasets to ensure models are evaluated using accurate assessments?

Errors in the InterCode Dataset by Type



Relationships between NL2SH datasets.



Dataset Findings

- **Research Question:** How can we validate NL2SH datasets to ensure models are evaluated using accurate assessments?
- **Result:** Manual verification of test data improves benchmark reliability. Cleaning and increasing training data improves in-weight learning.
- **Contribution:** We clean and increase the size of NL2SH test and training datasets by over 4x and 2x, respectively.

Benchmark Overview

- **Goal:**

Given			Determine
Natural Language	GT Command	Model Output	Equivalent
print the system boot time	who -b	uptime -s	True
delete bin\ in the current dir	rm -r ./bin	rm -r /bin	False
list all groups on the system	getent group	cat /etc/group	True

- **Challenge:** Benchmarks cannot determine the functional equivalence of Bash commands.
- **Research Question:** How can we design a functional equivalence heuristic that accurately measures the quality of model translations?

NL2SH Translation

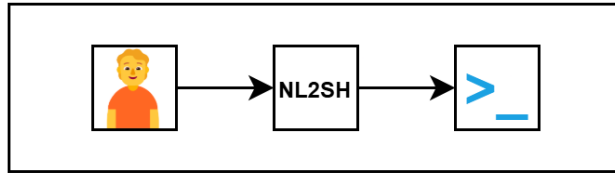


Diagram Key



Natural Language Prompt: Display the current time



NL2SH Model: Fine-tuned LLaMa3.1-8b-Instruct



NL2SH Model Command: timedatectl | grep "Local time"



NL2SH Model Command Output: 2024-06-27 14:18:48



Ground Truth Command: date +"%H:%M:%S"



Ground Truth Command Output: 14:18:48



Conventional Evaluation: BLEU, Edit, TF-IDF



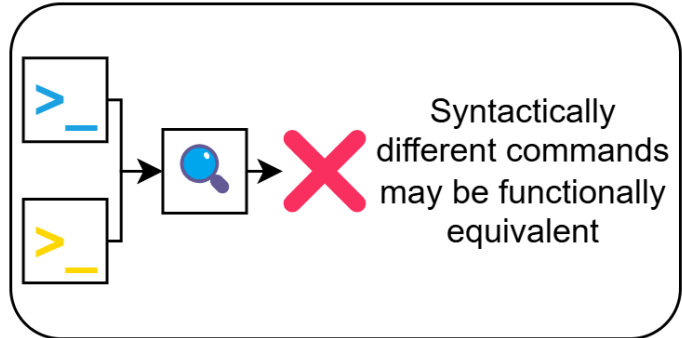
Command Execution Environment: Docker Container



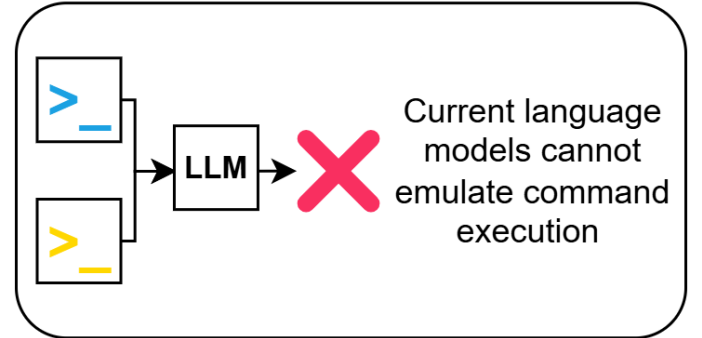
Language Model Evaluation: OpenAI's GPT-4

NL2SH Model Benchmarking Methods

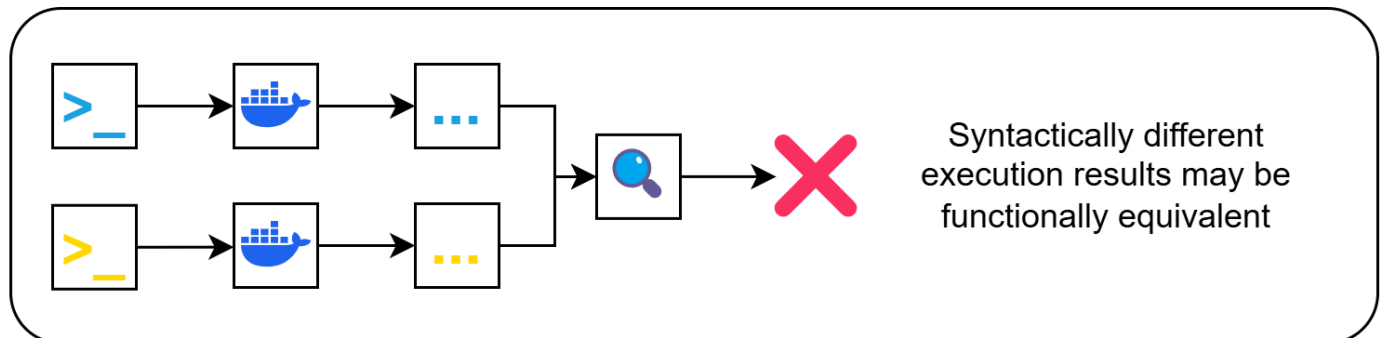
Conventional Evaluation (Agarwal et al.)



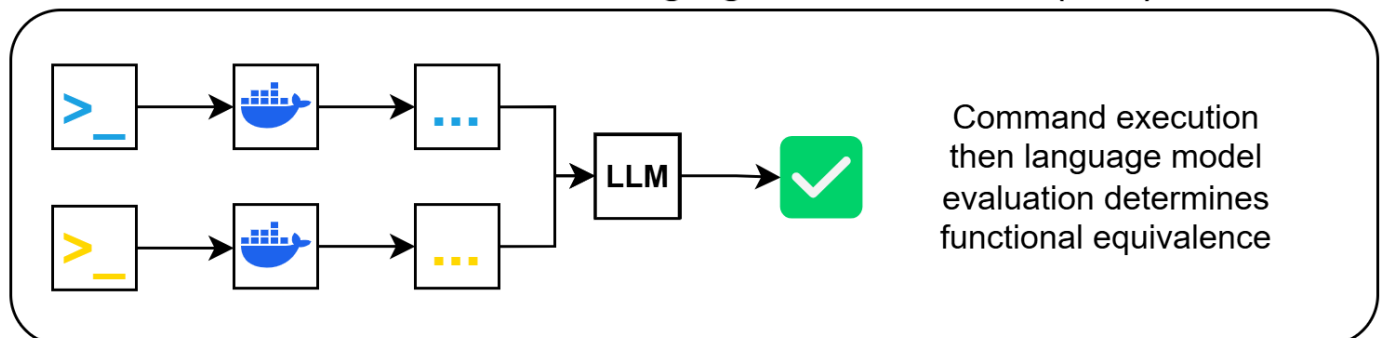
Language Model Evaluation (Song et al.)



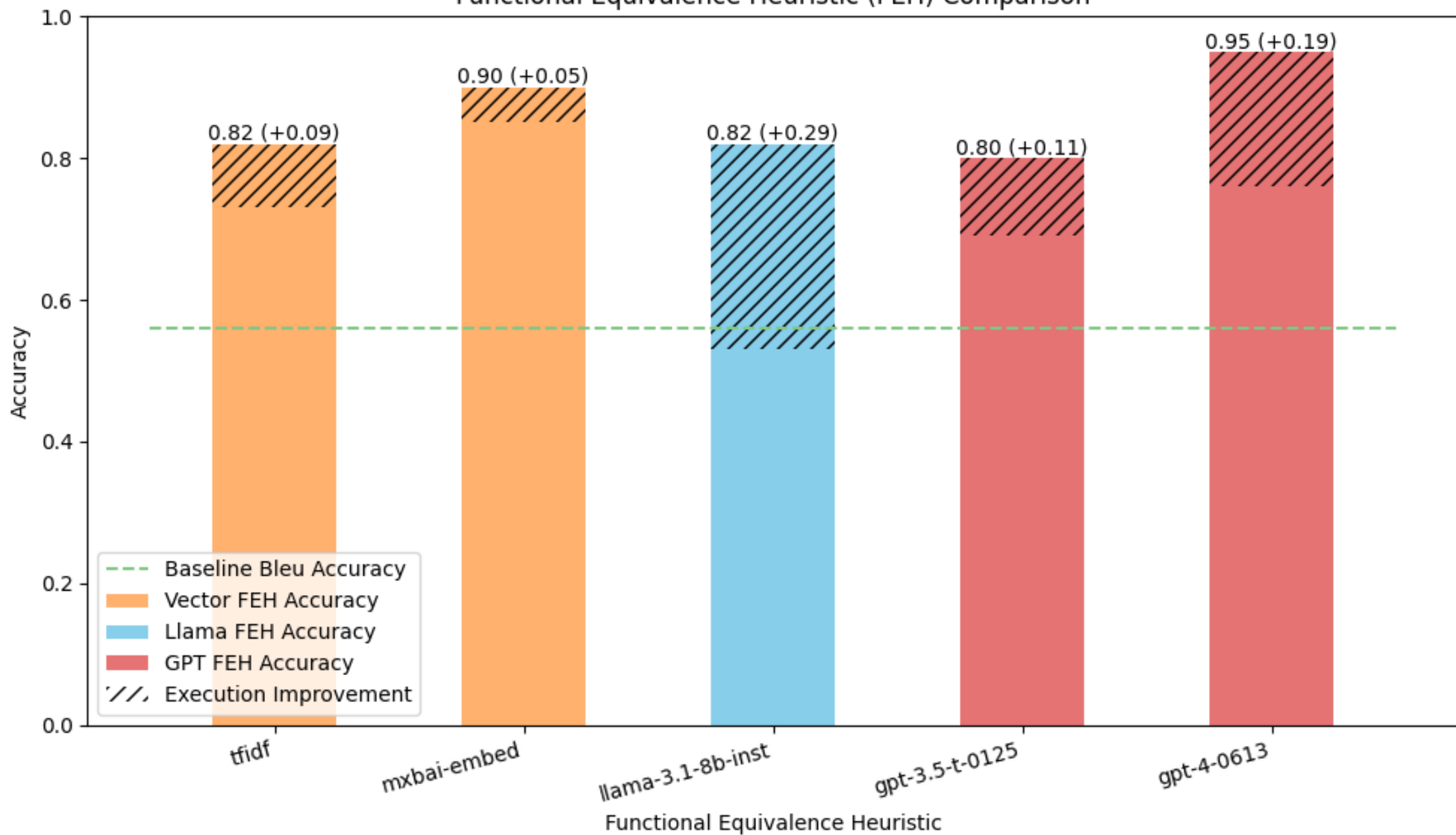
Command Execution + Conventional Evaluation (Yang et al.)



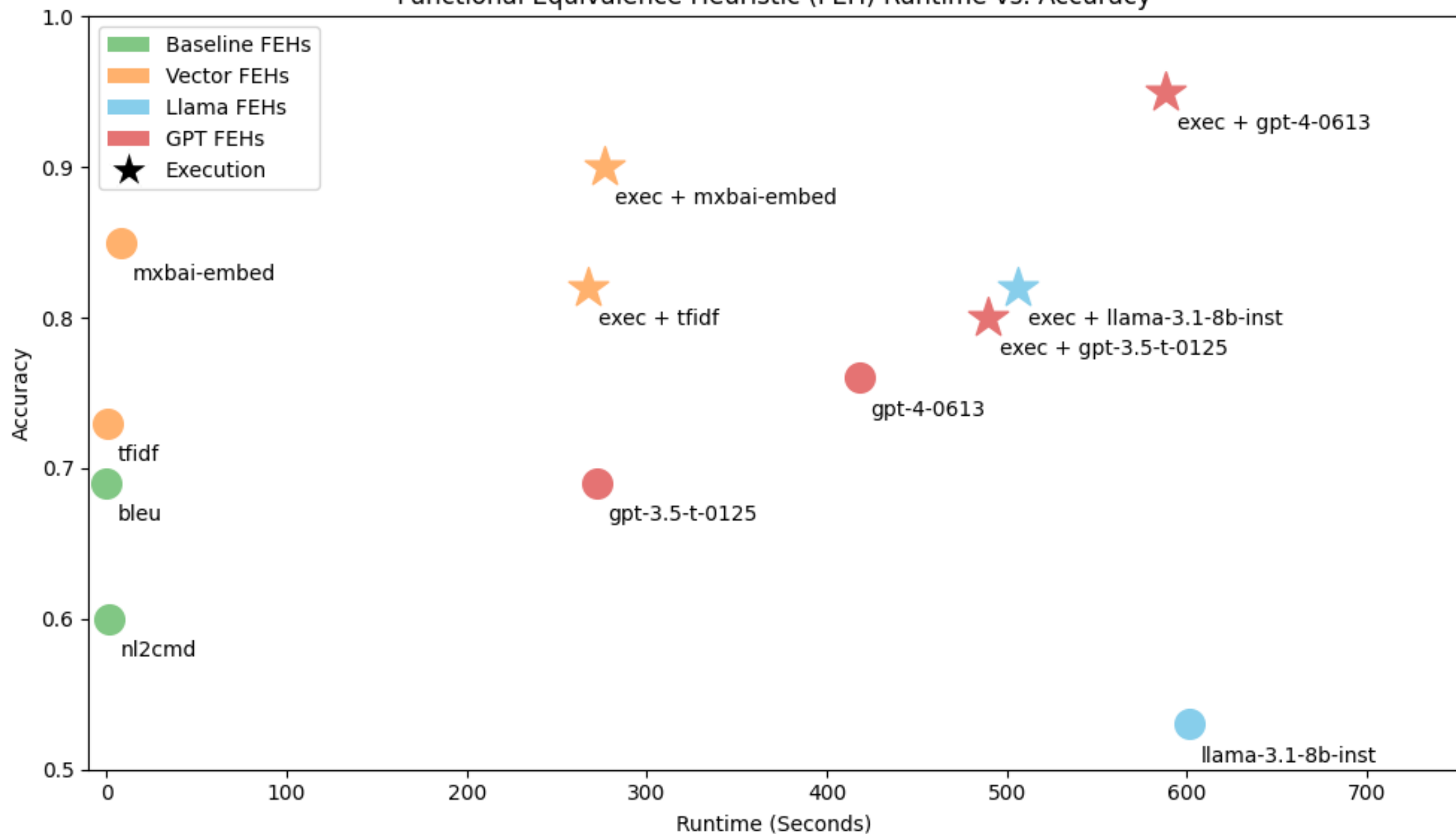
Command Execution + Language Model Evaluation (Ours)



Functional Equivalence Heuristic (FEH) Comparison



Functional Equivalence Heuristic (FEH) Runtime vs. Accuracy



Benchmark Findings

- **Research Question:** How can we design a functional equivalence heuristic that accurately measures the quality of model translations?
- **Result:** Command execution with language model evaluation of command outputs outperforms previous benchmarking methods.
- **Contribution:** We present a state-of-the-art method for determining the functional equivalence of Bash commands.

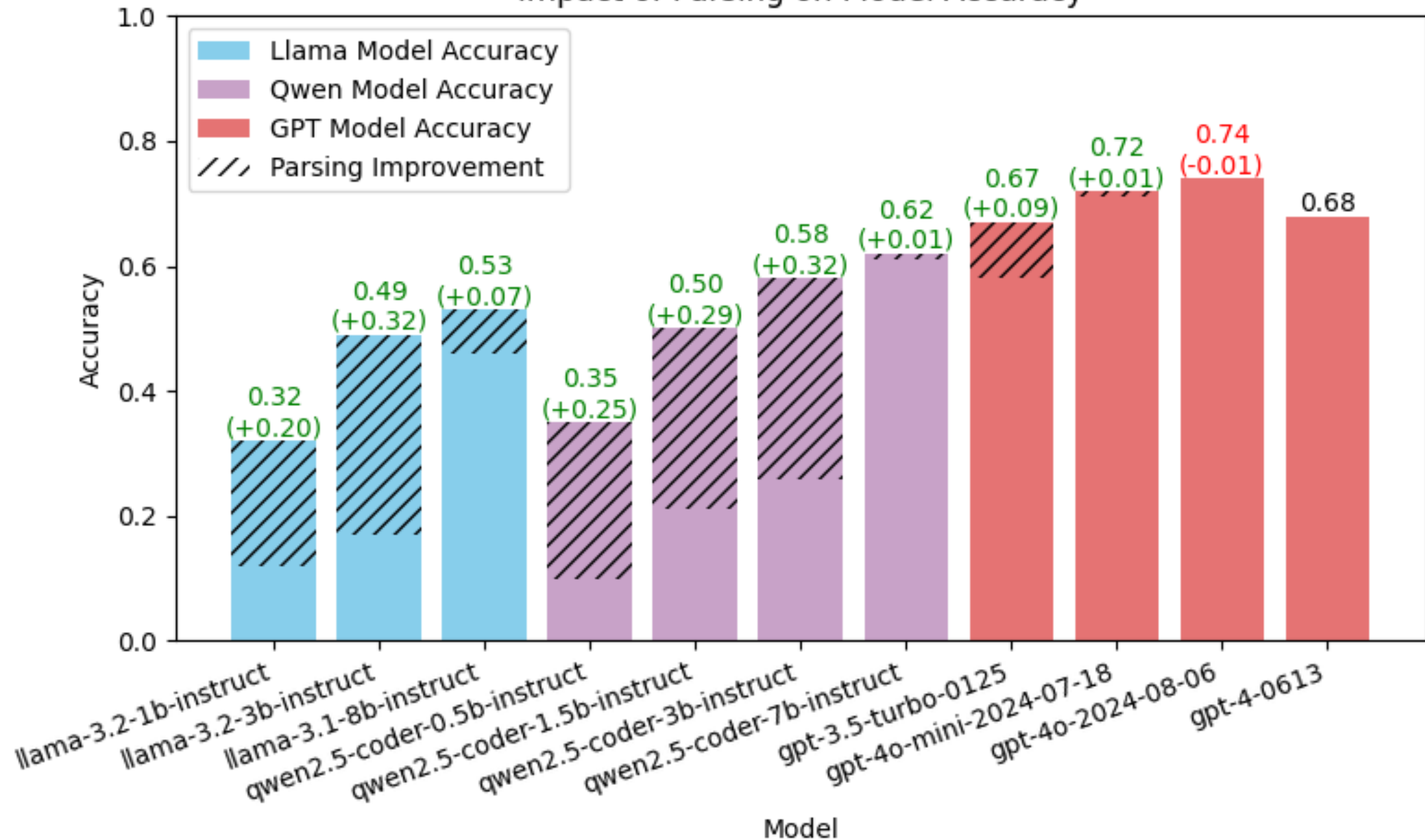
Translation Method Overview

- **Goal:**

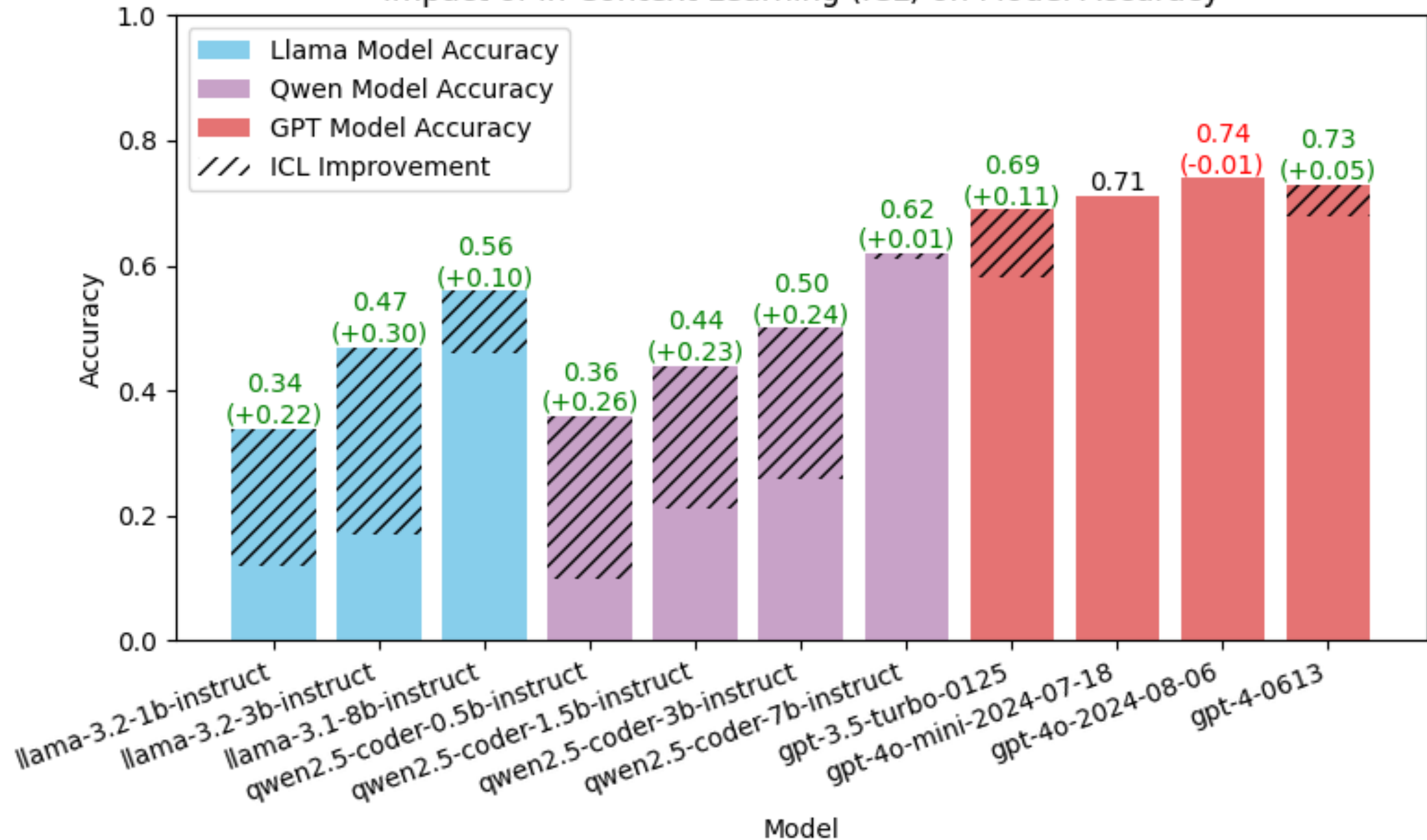
Given	Determine
Natural Language	Bash Command
List files in /workspace accessed over an hour ago base64 decode aGVsbG8 Sort and print only group names from /etc/group	<pre>find /workspace -type f -amin +60 echo 'aGVsbG8=' openssl enc -base64 -d cut -d: -f1 /etc/group sort</pre>

- **Challenge:** LLMs have poor NL2SH performance.
- **Research Question:** How can we improve the NL2SH performance of LLMs as measured by a reliable benchmark?

Impact of Parsing on Model Accuracy



Impact of In-Context Learning (ICL) on Model Accuracy



Translation Method Findings

- **Research Question:** How can we improve the accuracy of NL2SH models as measured by a reliable benchmark?
- **Result:** Parsing, in-context learning, in-weight learning and constrained decoding can improve NL2SH accuracy by up to 32%.
- **Contribution:** We demonstrate methods for improving the NL2SH performance of open and closed-source LLMs.

Conclusion

- Clean datasets are necessary for model training and evaluation
- Bash command execution paired with language model evaluation of command outputs can determine functional equivalence
- Parsing and in-context learning improve NL2SH performance
- NL2SH remains a difficult task and model outputs should be verified before they are used in real-world applications

Future Work

- Automate methods for dataset verification
- Evaluate the correctness of synthetic training data
- Increase speed and scalability of execution-based benchmarking
- Apply execution+LLM benchmarking to other translation tasks
- Evaluate the translation and functional equivalence capabilities of reasoning models
- Evaluate methods for increasing translation speed

Thank You!

