

Can LLMs translate SATisfactorily?

TL;DR

Assessing LLMs in Generating and Interpreting Formal Specifications

LLMs cannot translate formal syntax* but we can now automatically assess them! *yet

Rushang Karia, Daksh Dobhal, Daniel Bramblett, Pulkit Verma, Siddharth Srivastava

Motivation

- Use of LLMs to translate/interpret formal syntax is increasing
- McKinsey revised 50% automation estimates by a decade



Introducing **Devin**, the first AI software engineer

"Kids shouldn't learn to code"

- Jensen Huang, NVIDIA CEO

WORLD GOVERNMENT SUMMIT, DUBAI, 2024

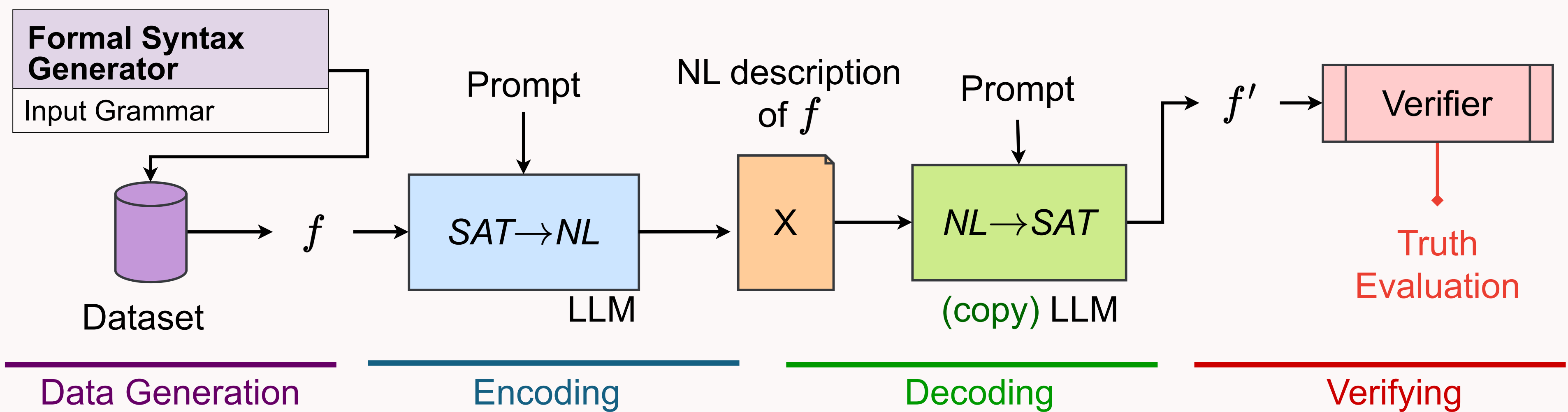
Widespread adoption: PRINCETON UNIVERSITY Microsoft ...

Key Challenges

How do we verify the translation capabilities of LLMs?

- How to efficiently generate new OOD data as LLMs evolve?
- How to annotate ground-truth data effectively?
- How to be robust in lieu of LLM hallucinations?

Can we make the pipeline automatic and human-independent?



Our Approach $NL \leftrightarrow SAT$

- Is scalable:** Uses syntax generators to scale datasets
- Is handsfree:** Uses two copies of an LLM to perform translation
- Is robust:** Uses external verifiers to validate the results

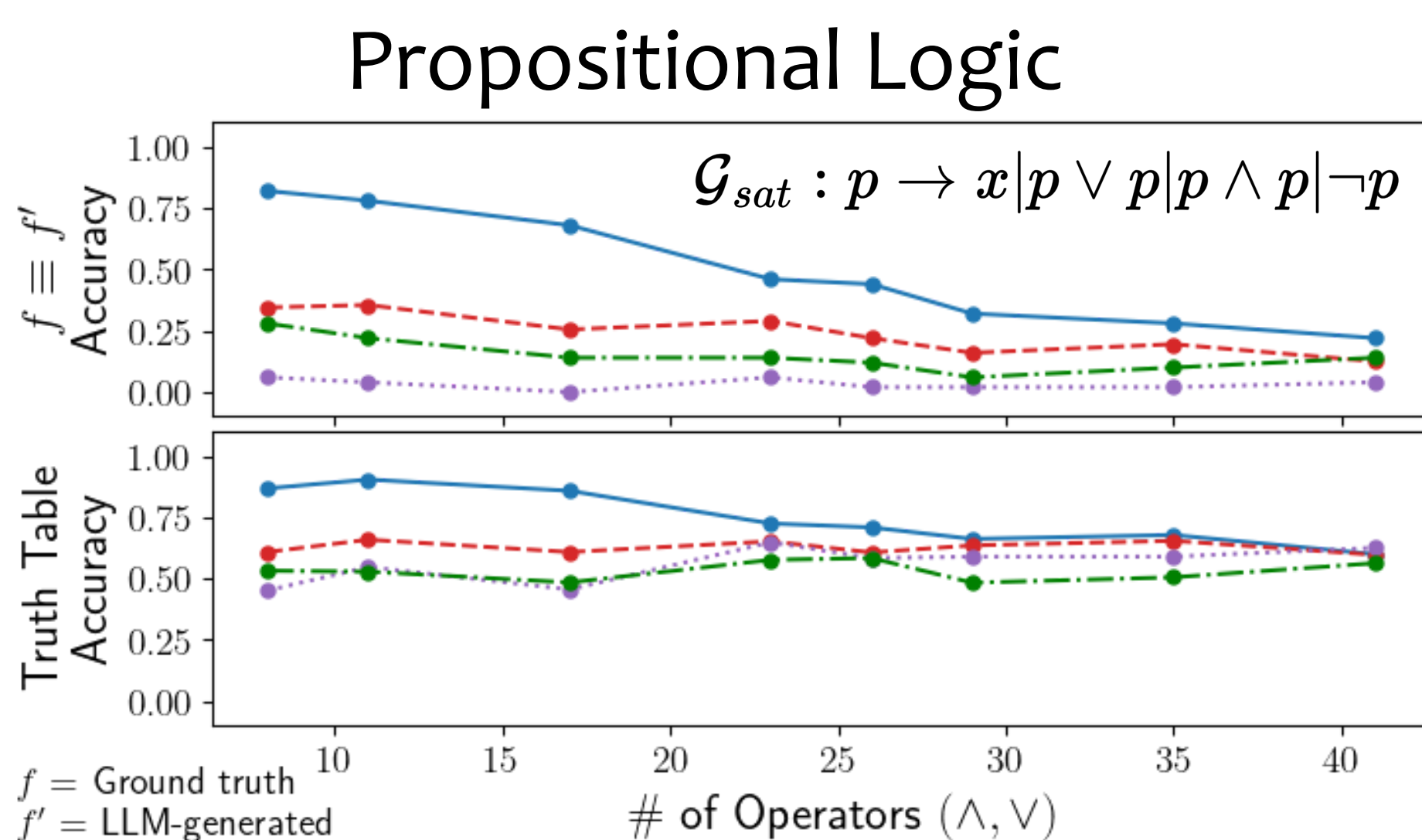
Higher values better

Legend:

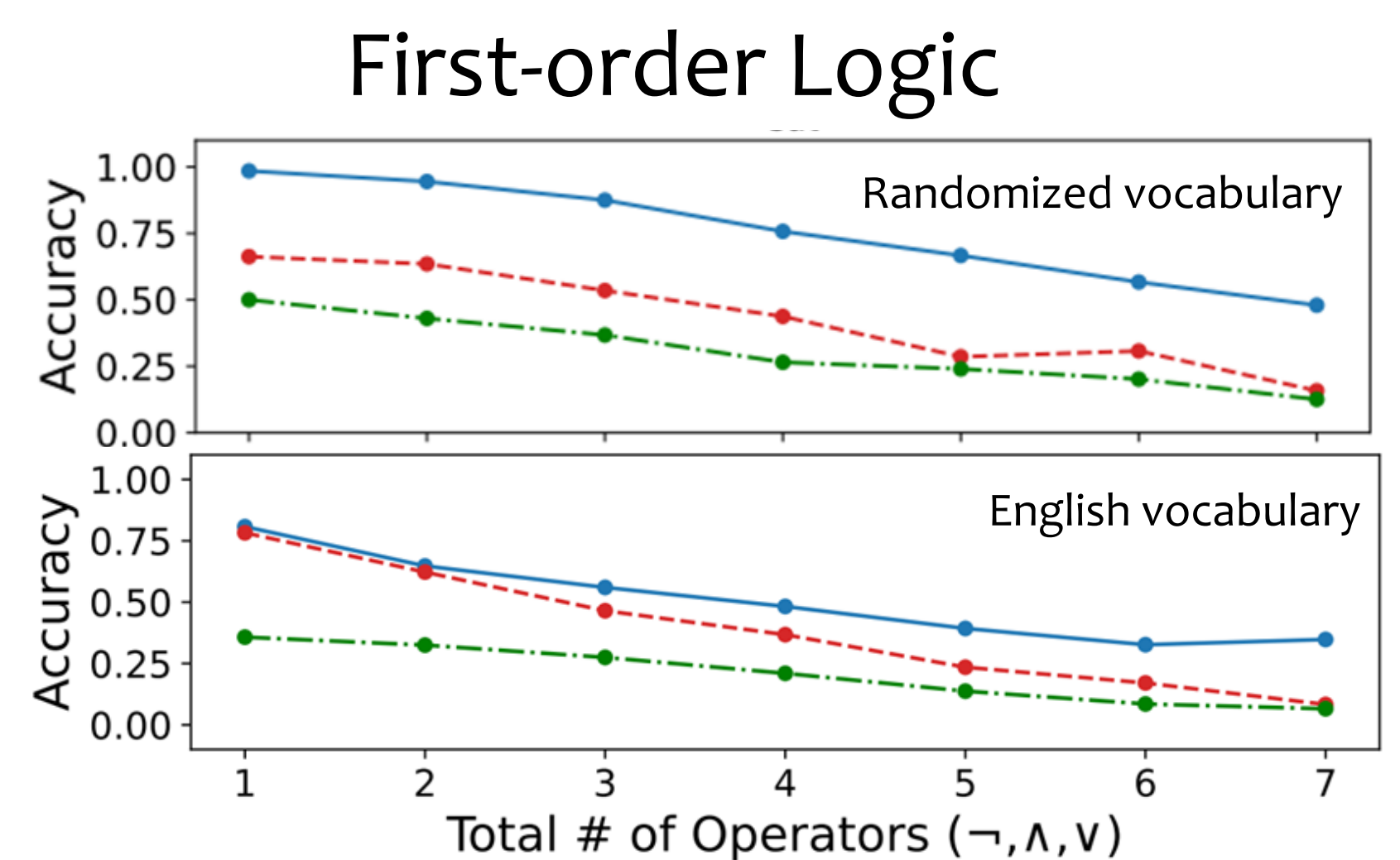
—●— GPT-4 —●— GPT-3.5-turbo —●— Mistral —●— Gemini

Prompt tuning: $\geq 95\%$ accuracy on k-SAT

Results



- Performance \downarrow as formula size \uparrow
- LLMs often misplace parentheses
- LLMs often hallucinate propositions



- Performance \downarrow as formula size \uparrow
- Worse than G_{sat} on smaller inputs
 - Similar failures as G_{sat}
 - Quantifiers often misplaced
- Using English is worse