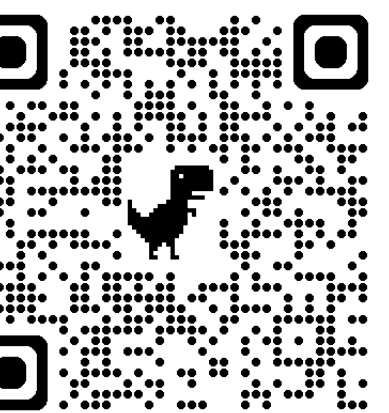


AnaloBench: Benchmarking the Identification of Abstract and Long-context Analogies



Do LLMs still perform well on challenging analogical reasoning tasks?

Motivation

AnaloBench moves beyond simple analogies for **challenging analogies** with **paragraphs** of raw-form text

Before - Simple Analogies

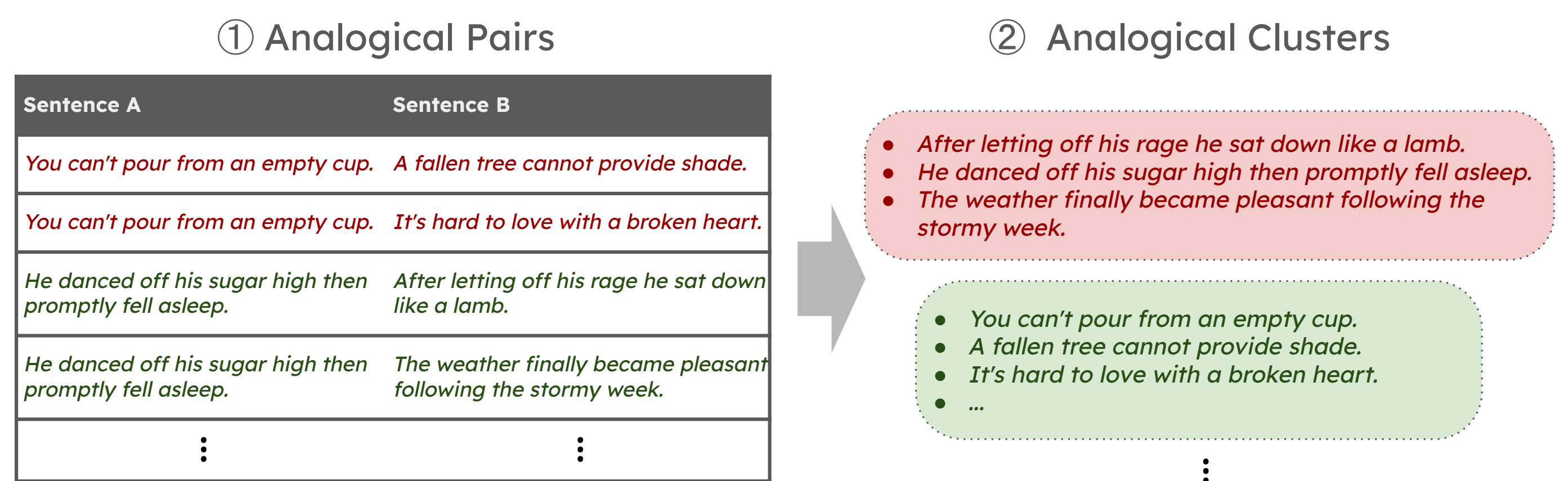
Text A: gust of wind : soft breeze Text B: fire : wisp of smoke

AnaloBench

Text A: In the heart of a small, poverty-stricken village, the weather began to take a frightful turn...The breeze gently kissed the cheeks of the shocked villagers, whispering an apology in their ears...
Text B: The candle had been burning for hours, casting a warm glow across the room... The smoke twirled and danced gracefully in the air, a subtle performance of nature, before slowly beginning to dissipate...

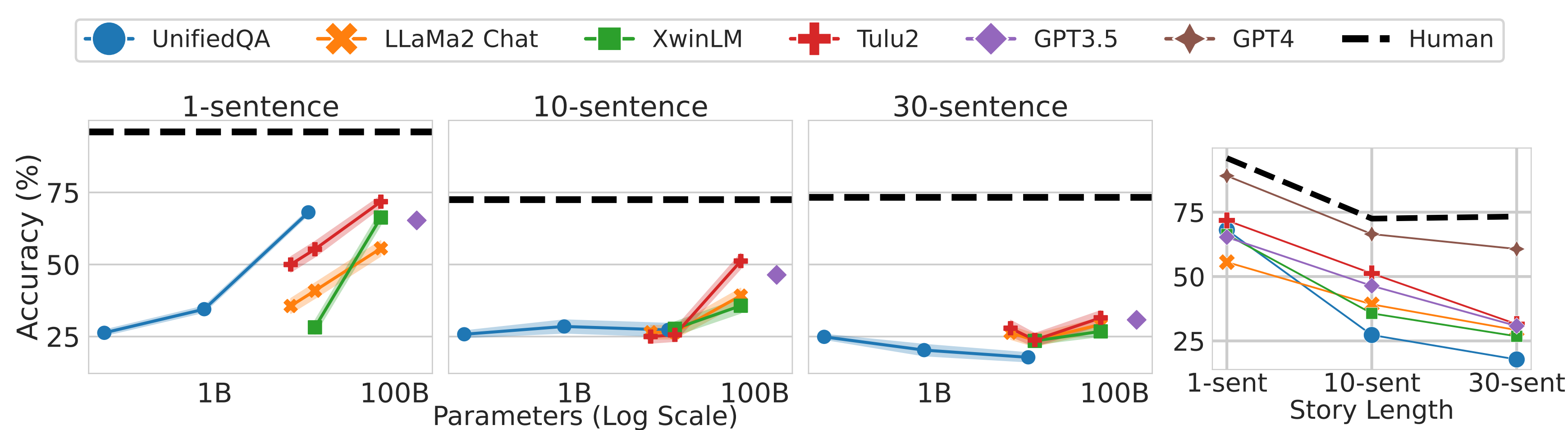
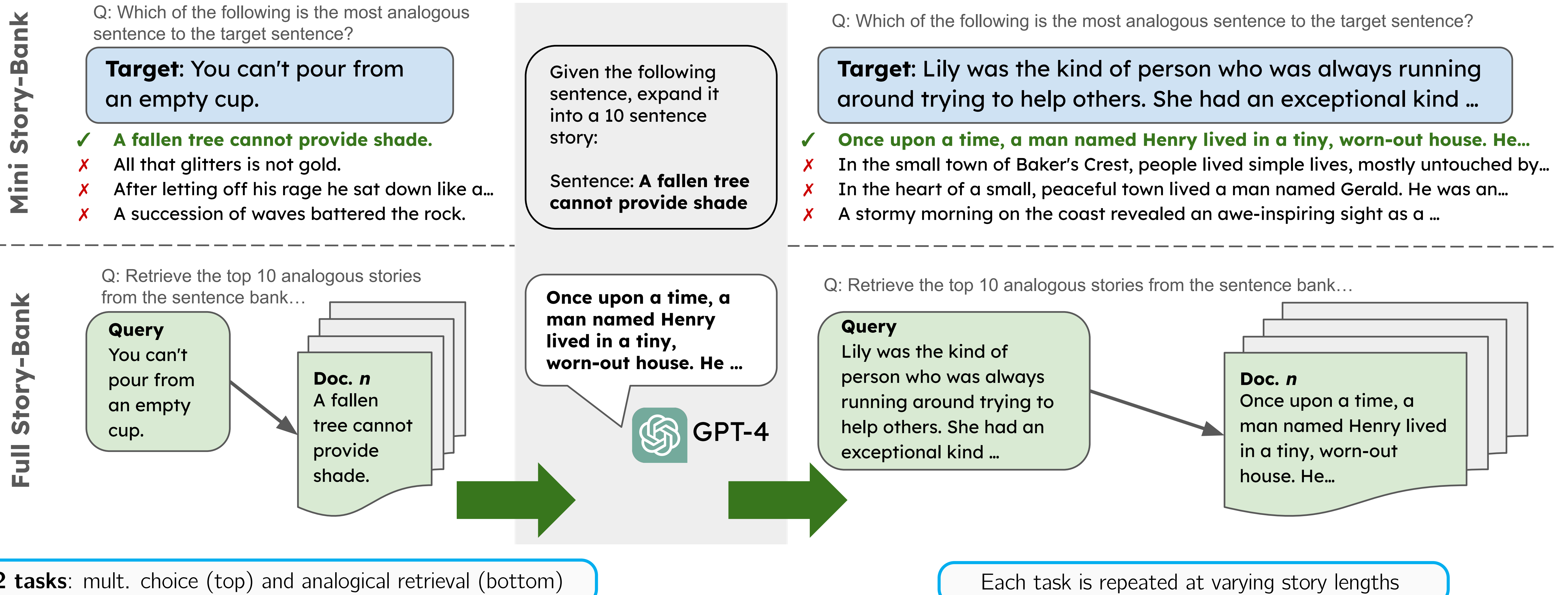
Dataset Creation

AnaloBench features 340 **human-contributed** analogies



- Each annotated analogy consists of a pair of analogous sentences
- Sentences are rewritten as 10 or 30 sentence narratives

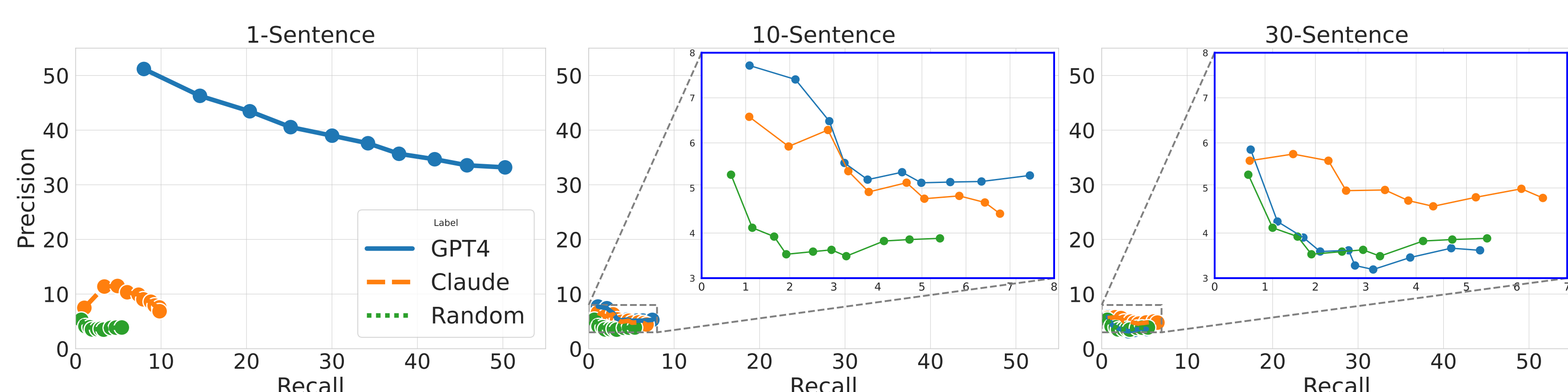
Benchmark Tasks



Model size is less helpful at longer narrative settings

↑ Human-AI gap

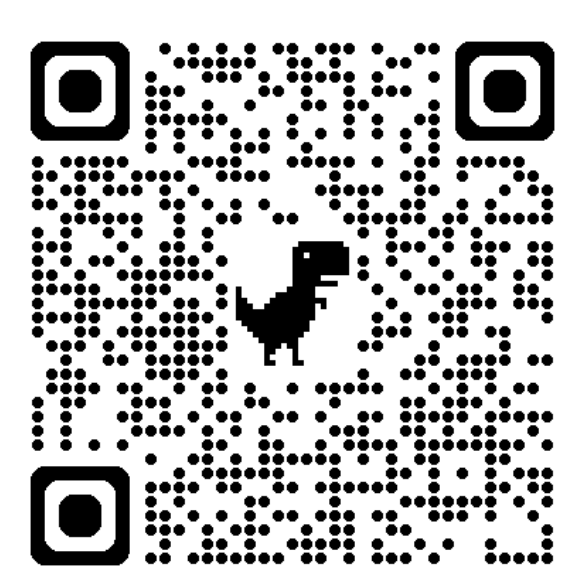
Analogy retrieval performance degrades with increased narrative length



Summary



Github



Dataset

Background

- Analogical reasoning was challenging for early AI systems
- Recent works claim that analogical reasoning in large language models is emergent behavior (Webb et al., 2023)

Findings

Do LLMs still perform well on challenging analogical reasoning tasks?

- Long narrative analogies are difficult for LLMs
- LLMs perform poorly on analogical retrieval