

GPT 4

GSM8K

MultispanQA

MMLU

50%

30%

10%

Hotpo

PIQA

RACE

Hellaswag

- Baseline understanding gap
- Prompt-improved understanding gap

