

	Claude 3.5 Sonnet	Claude 3 Opus	GPT-4o	Gemini 1.5 Pro
Graduate level reasoning <i>GPQA, Diamond</i>	<b>59.4%*</b> 0-shot CoT	<b>50.4%</b> 0-shot CoT	<b>53.6%</b> 0-shot CoT	—
Undergraduate level knowledge <i>MMLU</i>	<b>88.7%**</b> 5-shot	<b>86.8%</b> 5-shot	—	<b>85.9%</b> 5-shot
	<b>88.3%</b> 0-shot CoT	<b>85.7%</b> 0-shot CoT	<b>88.7%</b> 0-shot CoT	—
Code <i>HumanEval</i>	<b>92.0%</b> 0-shot	<b>84.9%</b> 0-shot	<b>90.2%</b> 0-shot	<b>84.1%</b> 0-shot
Multilingual math <i>MGSM</i>	<b>91.6%</b> 0-shot CoT	<b>90.7%</b> 0-shot CoT	<b>90.5%</b> 0-shot CoT	<b>87.5%</b> 8-shot
Reasoning over text <i>DROP, F1 score</i>	<b>87.1</b> 3-shot	<b>83.1</b> 3-shot	<b>83.4</b> 3-shot	<b>74.9</b> Variable shots
Mixed evaluations <i>BIG-Bench-Hard</i>	<b>93.1%</b> 3-shot CoT	<b>86.8%</b> 3-shot CoT	—	<b>89.2%</b> 3-shot CoT
Math problem-solving <i>MATH</i>	<b>71.1%</b> 0-shot CoT	<b>60.1%</b> 0-shot CoT	<b>76.6%</b> 0-shot CoT	<b>67.7%</b> 4-shot
Grade school math <i>GSM8K</i>	<b>96.4%</b> 0-shot CoT	<b>95.0%</b> 0-shot CoT	—	<b>90.8%</b> 11-shot

\* Claude 3.5 Sonnet scores 67.2% on 5-shot CoT GPQA with maj@32

\*\* Claude 3.5 Sonnet scores 90.4% on MMLU with 5-shot CoT prompting