

# Claude 3

Anthropic announces Claude 3, their new family of models that include Claude 3 Haiku, Claude 3 Sonnet, and Claude 3 Opus.

Claude 3 Opus (the strongest model) is reported to outperform GPT-4 and all other models on common benchmarks like MMLU and HumanEval.

## Results and Capabilities

---

Claude 3 capabilities include advanced reasoning, basic mathematics, analysis, data extraction, forecasting, content creation, code generation, and converting in non-English languages like Spanish, Japanese, and French. The table below demonstrates how Claude 3 compares with other models on several benchmarks with Claude 3 Opus outperforming all the mentioned models:

	Claude 3 Opus	Claude 3 Sonnet	Claude 3 Haiku	GPT-4	GPT-3.5	Gemini 1.0 Ultra	Gemini 1.0 Pro
Undergraduate level knowledge <i>MMLU</i>	<b>86.8%</b> 5-shot	<b>79.0%</b> 5-shot	<b>75.2%</b> 5-shot	<b>86.4%</b> 5-shot	<b>70.0%</b> 5-shot	<b>83.7%</b> 5-shot	<b>71.8%</b> 5-shot
Graduate level reasoning <i>GPQA, Diamond</i>	<b>50.4%</b> 0-shot CoT	<b>40.4%</b> 0-shot CoT	<b>33.3%</b> 0-shot CoT	<b>35.7%</b> 0-shot CoT	<b>28.1%</b> 0-shot CoT	—	—
Grade school math <i>GSM8K</i>	<b>95.0%</b> 0-shot CoT	<b>92.3%</b> 0-shot CoT	<b>88.9%</b> 0-shot CoT	<b>92.0%</b> 5-shot CoT	<b>57.1%</b> 5-shot	<b>94.4%</b> Maj1@32	<b>86.5%</b> Maj1@32
Math problem-solving <i>MATH</i>	<b>60.1%</b> 0-shot CoT	<b>43.1%</b> 0-shot CoT	<b>38.9%</b> 0-shot CoT	<b>52.9%</b> 4-shot	<b>34.1%</b> 4-shot	<b>53.2%</b> 4-shot	<b>32.6%</b> 4-shot
Multilingual math <i>MGSM</i>	<b>90.7%</b> 0-shot	<b>83.5%</b> 0-shot	<b>75.1%</b> 0-shot	<b>74.5%</b> 8-shot	—	<b>79.0%</b> 8-shot	<b>63.5%</b> 8-shot
Code <i>HumanEval</i>	<b>84.9%</b> 0-shot	<b>73.0%</b> 0-shot	<b>75.9%</b> 0-shot	<b>67.0%</b> 0-shot	<b>48.1%</b> 0-shot	<b>74.4%</b> 0-shot	<b>67.7%</b> 0-shot
Reasoning over text <i>DROP, F1 score</i>	<b>83.1</b> 3-shot	<b>78.9</b> 3-shot	<b>78.4</b> 3-shot	<b>80.9</b> 3-shot	<b>64.1</b> 3-shot	<b>82.4</b> Variable shots	<b>74.1</b> Variable shots
Mixed evaluations <i>BIG-Bench-Hard</i>	<b>86.8%</b> 3-shot CoT	<b>82.9%</b> 3-shot CoT	<b>73.7%</b> 3-shot CoT	<b>83.1%</b> 3-shot CoT	<b>66.6%</b> 3-shot CoT	<b>83.6%</b> 3-shot CoT	<b>75.0%</b> 3-shot CoT
Knowledge Q&A <i>ARC-Challenge</i>	<b>96.4%</b> 25-shot	<b>93.2%</b> 25-shot	<b>89.2%</b> 25-shot	<b>96.3%</b> 25-shot	<b>85.2%</b> 25-shot	—	—
Common Knowledge <i>HellaSwag</i>	<b>95.4%</b> 10-shot	<b>89.0%</b> 10-shot	<b>85.9%</b> 10-shot	<b>95.3%</b> 10-shot	<b>85.5%</b> 10-shot	<b>87.8%</b> 10-shot	<b>84.7%</b> 10-shot

Claude 3 Haiku is the fastest and most cost-effective model of the series. Claude 3 Sonnet is 2x faster than previous iterations of Claude and Opus is as fast as Claude 2.1 with more superior capabilities.

The Claude 3 models offer support for 200K context windows but can be extended to 1M tokens to select customers. Claude 3 Opus achieved near-perfect recall on the Needle In A Haystack (NIAH) evaluation which measures the model's ability to recall information in a large corpus and effectively process long context prompts.

The models also have strong vision capabilities for processing formats like photos, charts, and graphs.

	Claude 3 Opus	Claude 3 Sonnet	Claude 3 Haiku	GPT-4V	Gemini 1.0 Ultra	Gemini 1.0 Pro
Math & reasoning <i>MMMU (val)</i>	59.4%	53.1%	50.2%	56.8%	59.4%	47.9%
Document visual Q&A <i>ANLS score, test</i>	89.3%	89.5%	88.8%	88.4%	90.9%	88.1%
Math <i>MathVista (testmini)</i>	50.5% CoT	47.9% CoT	46.4% CoT	49.9%	53.0%	45.2%
Science diagrams <i>AI2D, test</i>	88.1%	88.7%	86.7%	78.2%	79.5%	73.9%
Chart Q&A <i>Relaxed accuracy (test)</i>	80.8% 0-shot CoT	81.1% 0-shot CoT	81.7% 0-shot CoT	78.5% 4-shot CoT	80.8%	74.1%

Anthropic also claim that these models have a more nuanced understanding of requests and make fewer refusals. Opus also shows significant improvements in factual question answering in open-ended questions while reducing incorrect answers or hallucinations. Claude 3 models are also better than the Claude 2 models at producing structured outputs like JSON objects.

## References

- [Claude 3 Haiku, Claude 3 Sonnet, and Claude 3 Opus](#)
- [The Claude 3 Model Family: Opus, Sonnet, Haiku](#)

Last updated on September 19, 2024