

# Gemini: A Family of Highly Capable Multimodal Models

Gemini Team, Google<sup>1</sup>

This report introduces a new family of multimodal models, Gemini, that exhibit remarkable capabilities across image, audio, video, and text understanding. The Gemini family consists of Ultra, Pro, and Nano sizes, suitable for applications ranging from complex reasoning tasks to on-device memory-constrained use-cases. Evaluation on a broad range of benchmarks shows that our most-capable Gemini Ultra model advances the state of the art in 30 of 32 of these benchmarks — notably being the first model to achieve human-expert performance on the well-studied exam benchmark MMLU, and improving the state of the art in every one of the 20 multimodal benchmarks we examined. We believe that the new capabilities of the Gemini family in cross-modal reasoning and language understanding will enable a wide variety of use cases. We discuss our approach toward post-training and deploying Gemini models responsibly to users through services including Gemini, Gemini Advanced, Google AI Studio, and Cloud Vertex AI.

## 1. Introduction

We present Gemini, a family of highly capable multimodal models developed at Google. We trained Gemini models jointly across image, audio, video, and text data for the purpose of building a model with both strong generalist capabilities across modalities alongside cutting-edge understanding and reasoning performance in each respective domain.

Gemini 1.0, our first version, comes in three sizes: Ultra for highly-complex tasks, Pro for enhanced performance and deployability at scale, and Nano for on-device applications. Each size is specifically tailored to address different computational limitations and application requirements.

After large-scale pre-training, we post-train our models to improve overall quality, enhance target capabilities, and ensure alignment and safety criteria are met. Due to the varied requirements of our downstream applications, we have produced two post-trained Gemini model family variants. Chat-focused variants, referred to as Gemini Apps models, are optimized for [Gemini and Gemini Advanced](#), our conversational AI service formerly known as Bard. Developer-focused variants, referred to as Gemini API models, are optimized for a range of products and are accessible through [Google AI Studio](#) and [Cloud Vertex AI](#).

We evaluate the performance of pre- and post-trained Gemini models on a comprehensive suite of internal and external benchmarks covering a wide range of language, coding, reasoning, and multimodal tasks.

The Gemini family advances state-of-the-art in large-scale language modeling ([Anil et al., 2023](#); [Brown et al., 2020](#); [Chowdhery et al., 2023](#); [Hoffmann et al., 2022](#); [OpenAI, 2023a](#); [Radford et al., 2019](#); [Rae et al., 2021](#)), image understanding ([Alayrac et al., 2022](#); [Chen et al., 2022](#); [Dosovitskiy et al., 2020](#); [OpenAI, 2023b](#); [Reed et al., 2022](#); [Yu et al., 2022a](#)), audio processing ([Radford et al., 2023](#); [Zhang et al., 2023](#)), and video understanding ([Alayrac et al., 2022](#); [Chen et al., 2023](#)). It also builds on the work on sequence models ([Sutskever et al., 2014](#)), a long history of work in deep learning based on neural networks ([LeCun et al., 2015](#)), and machine learning distributed systems

---

<sup>1</sup>See Contributions and Acknowledgments section for full author list. Please send correspondence to [gemini-1-report@google.com](mailto:gemini-1-report@google.com)

(Barham et al., 2022; Bradbury et al., 2018; Dean et al., 2012) that enable large-scale training.

Our most capable model, Gemini Ultra, achieves new state-of-the-art results in 30 of 32 benchmarks we report on, including 10 of 12 popular text and reasoning benchmarks, 9 of 9 image understanding benchmarks, 6 of 6 video understanding benchmarks, and 5 of 5 speech recognition and speech translation benchmarks. Gemini Ultra is the first model to achieve human-expert performance on MMILU (Hendrycks et al., 2021a) — a prominent benchmark testing knowledge and reasoning via a suite of exams — with a score above 90%. Beyond text, Gemini Ultra makes notable advances on challenging multimodal reasoning tasks. For example, on the recent MMMU benchmark (Yue et al., 2023), that comprises questions about images on multi-discipline tasks requiring college-level subject knowledge and deliberate reasoning, Gemini Ultra achieves a new state-of-the-art score of 62.4%, outperforming the previous best model by more than 5 percentage points. It provides a uniform performance lift for video question answering and audio understanding benchmarks.

Qualitative evaluation showcases impressive crossmodal reasoning capabilities, enabling the model to understand and reason across an input sequence of audio, images, and text natively (see Figure 5 and Table 13). Consider the educational setting depicted in Figure 1 as an example. A teacher has drawn a physics problem of a skier going down a slope, and a student has worked through a solution to it. Using Gemini models’ multimodal reasoning capabilities, the model is able to understand the messy handwriting, correctly understand the problem formulation, convert both the problem and solution to mathematical typesetting, identify the specific step of reasoning where the student went wrong in solving the problem, and then give a worked through correct solution to the problem. This opens up exciting educational possibilities, and we believe the new multimodal and reasoning capabilities of Gemini models have dramatic applications across many fields.

The reasoning capabilities of large language models show promise toward building generalist agents that can tackle more complex multi-step problems. The AlphaCode team built AlphaCode 2 (Leblond et al, 2023), a new Gemini-model-powered agent, that combines Gemini models’ reasoning capabilities with search and tool-use to excel at solving competitive programming problems. AlphaCode 2 ranks within the top 15% of entrants on the Codeforces competitive programming platform, a large improvement over its state-of-the-art predecessor in the top 50% (Li et al., 2022).

In tandem, we advance the frontier of efficiency with Gemini Nano, a series of small models targeting on-device deployment. These models excel in on-device tasks, such as summarization, reading comprehension, text completion tasks, and exhibit impressive capabilities in reasoning, STEM, coding, multimodal, and multilingual tasks relative to their sizes.

In the following sections, we first provide an overview of the model architecture, training infrastructure, and pre-training dataset. We then present detailed *evaluations* of the pre- and post-trained Gemini model family, covering well-studied benchmarks across text, code, image, audio and video — which include both English performance and multilingual capabilities. Next we discuss our approach to post-training, highlight common and distinct aspects of the Gemini Apps and Gemini API model variants, and benchmark their performance on key capabilities. *Responsible deployment* is critical: we explain our process for impact assessments, developing model policies, evaluations, and mitigations of harm before deployment decisions. Finally, we discuss the broader implications of Gemini models, their limitations alongside their potential applications — paving the way for a new era of research and innovation in AI.

**Prompt**

Here's a solution to a Physics problem by a student.

23. A skier slides down the frictionless slope as shown. What is the skier's speed at the bottom?

Solution:

The conservation of energy means that the energy of the skier is same at the beginning and at the end.

Start:  $E = mgL$

End:  $E = \frac{1}{2}mv^2$

$$\Rightarrow mgL = \frac{1}{2}mv^2$$

$$\Rightarrow v = \sqrt{2gL} = \sqrt{2 \times 9.81 \times 80} = 39.6 \text{ m.s}^{-1}$$

Try to reason about the question step by step.

Did the student get the correct answer? If the solution is wrong, please explain what is wrong and solve the problem.

Make sure to use LaTeX for math and round off the final answer to two decimal places.

**Response**

The student did not get the correct answer. The student made a mistake in the calculation of the potential energy at the start of the slope. The potential energy at the start of the slope is given by  $mgh$ , where  $m$  is the mass of the skier,  $g$  is the acceleration due to gravity, and  $h$  is the height of the slope. The student used the length of the slope,  $L$ , instead of the height,  $H$ , in the calculation of the potential energy.

The correct solution is:

The conservation of energy means that the total energy of the skier is the same at the beginning and at the end of the slope.

At the start of the slope, the skier has only potential energy, which is given by  $mgh$ .

At the end of the slope, the skier has only kinetic energy, which is given by  $\frac{1}{2}mv^2$ .

Therefore, we can write:

$$mgh = \frac{1}{2}mv^2$$

Cancelling out the mass  $m$  on both sides, we get:

$$gh = \frac{1}{2}v^2$$

Solving for  $v$ , we get:

$$v = \sqrt{2gh}$$

Substituting the given values, we get:

$$v = \sqrt{2 \times 9.81 \times 40} = 28.01 \text{ m/s}$$

Therefore, the skier's speed at the bottom of the slope is 28.01 m/s.

Figure 1 | Verifying a student’s solution to a physics problem. The model is able to correctly recognize all of the handwritten content and verify the reasoning. On top of understanding the text in the image, it needs to understand the problem setup and correctly follow instructions to generate  $\text{\LaTeX}$ .

## 2. Model Architecture

Gemini models build on top of Transformer decoders (Vaswani et al., 2017b) that are enhanced with improvements in architecture and model optimization to enable stable training at scale and optimized inference on Google’s Tensor Processing Units. They are trained to support 32k context length, employing efficient attention mechanisms (for e.g. multi-query attention (Shazeer, 2019a)). Our first version, Gemini 1.0, comprises three main sizes to support a wide range of applications as discussed in Table 1.

Gemini models are trained to accommodate textual input interleaved with a wide variety of audio and visual inputs, such as natural images, charts, screenshots, PDFs, and videos, and they can produce text and image outputs (see Figure 2). The visual encoding of Gemini models is inspired by our own foundational work on Flamingo (Alayrac et al., 2022), CoCa (Yu et al., 2022a), and PaLI (Chen et al., 2022), with the important distinction that the models are multimodal from the beginning and can natively output images using discrete image tokens (Ramesh et al., 2021; Yu et al., 2022b).

Video understanding is accomplished by encoding the video as a sequence of frames in the large context window. Video frames or images can be interleaved naturally with text or audio as part of the model input. The models can handle variable input resolution in order to spend more compute on tasks that require fine-grained understanding. In addition, Gemini models can directly ingest audio

Model size	Model description
Ultra	Our most capable model that delivers state-of-the-art performance across a wide range of highly complex tasks, including reasoning and multimodal tasks. It is efficiently serveable at scale on TPU accelerators due to the Gemini architecture.
Pro	A performance-optimized model in terms of cost as well as latency that delivers significant performance across a wide range of tasks. This model exhibits strong reasoning performance and broad multimodal capabilities.
Nano	Our most efficient model, designed to run on-device. We trained two versions of Nano, with 1.8B (Nano-1) and 3.25B (Nano-2) parameters, targeting low and high memory devices respectively. It is trained by distilling from larger Gemini models. It is 4-bit quantized for deployment and provides best-in-class performance.

Table 1 | An overview of the Gemini 1.0 model family.

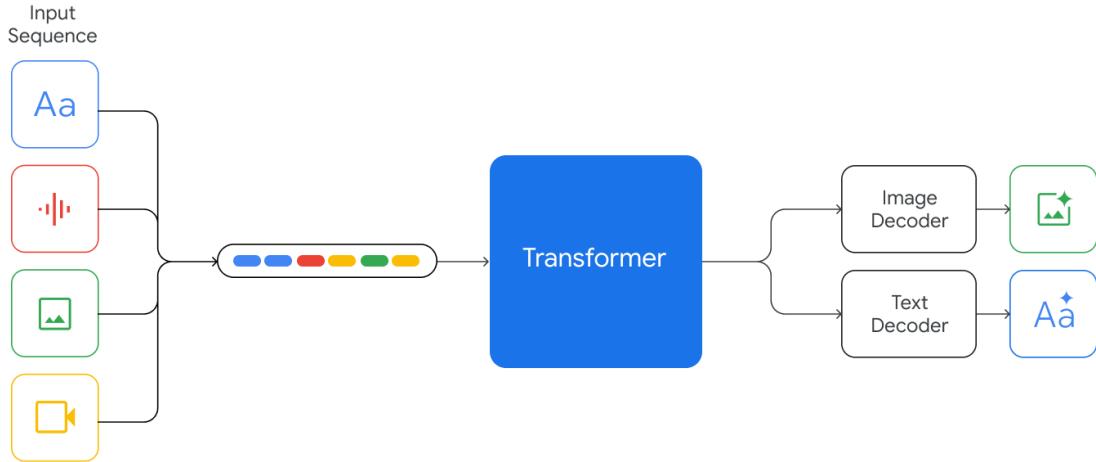


Figure 2 | Gemini models support interleaved sequences of text, image, audio, and video as inputs (illustrated by tokens of different colors in the input sequence). They can output responses with interleaved image and text.

signals at 16kHz from Universal Speech Model (USM) ([Zhang et al., 2023](#)) features. This enables the model to capture nuances that are typically lost when the audio is naively mapped to a text input (for example, see audio understanding demo on the [website](#)).

Training the Gemini family of models required innovations in training algorithms, dataset, and infrastructure. For the Pro model, the inherent scalability of our infrastructure and learning algorithms enable us to complete pre-training in a matter of weeks, leveraging a fraction of the Ultra’s resources. The Nano series of models leverage additional advancements in distillation and training algorithms to produce the best-in-class small language models for a wide variety of tasks, such as summarization and reading comprehension, which power our next generation on-device experiences.

### 3. Training Infrastructure

We trained Gemini models using TPUv5e and TPUv4 ([Jouppi et al., 2023](#)), depending on their sizes and configuration. Training Gemini Ultra used a large fleet of TPUv4 accelerators owned by Google

across multiple datacenters. This represents a significant increase in scale over our prior flagship model PaLM-2 which presented new infrastructure challenges. Scaling up the number of accelerators results in a proportionate decrease in the mean time between failure of hardware in the overall system. We minimized the rate of planned reschedules and preemptions, but genuine machine failures are commonplace across all hardware accelerators at such large scales.

TPUv4 accelerators are deployed in “SuperPods” of 4096 chips, each connected to a dedicated optical switch, which can dynamically reconfigure 4x4x4 chip cubes into arbitrary 3D torus topologies in around 10 seconds (Jouppi et al., 2023). For Gemini Ultra, we decided to retain a small number of cubes per superpod to allow for hot standbys and rolling maintenance.

TPU accelerators primarily communicate over the high speed inter-chip-interconnect, but at Gemini Ultra scale, we combine SuperPods in multiple datacenters using Google’s intra-cluster and inter-cluster network (Poutievski et al., 2022; Wetherall et al., 2023; Yao Hong et al., 2018). Google’s network latencies and bandwidths are sufficient to support the commonly used synchronous training paradigm, exploiting model parallelism within superpods and data-parallelism across superpods.

The ‘single controller’ programming model of Jax (Bradbury et al., 2018) and Pathways (Barham et al., 2022) allows a single Python process to orchestrate the entire training run, dramatically simplifying the development workflow. The GSPMD partitioner (Xu et al., 2021) in the XLA compiler partitions the training step computation, and the MegaScale XLA compiler (XLA, 2019) pass statically schedules appropriate collectives so that they maximally overlap with the computation with very little variation in step time.

Maintaining a high goodput<sup>2</sup> at this scale would have been impossible using the conventional approach of periodic checkpointing of weights to persistent cluster storage. For Gemini models, we instead made use of redundant in-memory copies of the model state, and on any unplanned hardware failures, we rapidly recover directly from an intact model replica. Compared to both PaLM and PaLM-2 (Anil et al., 2023), this provided a substantial speedup in recovery time, despite the significantly larger training resources being used. As a result, the overall goodput for the largest-scale training job increased from 85% to 97%.

Training at unprecedented scale invariably surfaces new and interesting systems failure modes - and in this instance one of the problems that we needed to address was that of “Silent Data Corruption (SDC)” (Dixit et al., 2021; Hochschild et al., 2021; Vishwanathan et al., 2015). Although these are extremely rare, the scale of Gemini models means that we can expect SDC events to impact training every week or two. Rapidly detecting and removing faulty hardware required several new techniques that exploit deterministic replay to isolate incorrect computations, combined with proactive SDC scanners on idle machines and hot standbys. Our fully deterministic infrastructure allowed us to quickly identify root causes (including hardware failures) during the development leading up to the Ultra model, and this was a crucial ingredient towards stable training.

## 4. Pre-Training Dataset

Gemini models are trained on a dataset that is both multimodal and multilingual. Our pre-training dataset uses data from web documents, books, and code, and includes image, audio, and video data.

We use the SentencePiece tokenizer (Kudo and Richardson, 2018) and find that training the tokenizer on a large sample of the entire training corpus improves the inferred vocabulary and subsequently improves model performance. For example, we find Gemini models can efficiently

---

<sup>2</sup>We define goodput as the time spent computing useful new steps over the elapsed time of the training job.

tokenize non-Latin scripts which can, in turn, benefit model quality as well as training and inference speed.

The number of tokens used to train the largest models were determined following the approach in [Hoffmann et al. \(2022\)](#). The smaller models are trained for significantly more tokens to improve performance for a given inference budget, similar to the approach advocated in [Touvron et al. \(2023a\)](#).

We apply quality filters to all datasets, using both heuristic rules and model-based classifiers. We also perform safety filtering to remove harmful content based on our policies. To maintain the integrity of evaluations, we search for and remove any evaluation data that may have been in our training corpus before using data for training. The final data mixtures and weights were determined through ablations on smaller models. We stage training to alter the mixture composition during training – increasing the weight of domain-relevant data towards the end of training. We find that data quality is an important factor for highly-performing models, and believe that many interesting questions remain around finding the optimal dataset distribution for pre-training.

## 5. Evaluation

The Gemini models are natively multimodal, as they are trained jointly across text, image, audio, and video. One open question is whether this joint training can result in a model which has strong capabilities in each domain – even when compared to models and approaches that are narrowly tailored to single domains. We find this to be the case: Gemini models set a new state of the art across a wide range of text, image, audio, and video benchmarks. [ww](#)

### 5.1. Text

#### 5.1.1. Academic Benchmarks

We compare pre- and post-trained Gemini Pro and Ultra models to a suite of external LLMs and our previous best model PaLM 2 across a series of text-based academic benchmarks covering reasoning, reading comprehension, STEM, and coding. We report these results in Table 2. Broadly, we find that the performance of Gemini Pro outperforms inference-optimized models such as GPT-3.5 and performs comparably with several of the most capable models available, and Gemini Ultra outperforms all current models. In this section, we examine some of these findings.

On MMLU ([Hendrycks et al., 2021a](#)), Gemini Ultra can outperform all existing models, achieving an accuracy of 90.04%. MMLU is a holistic exam benchmark, which measures knowledge across a set of 57 subjects. Human expert performance is gauged at 89.8% by the benchmark authors, and Gemini Ultra is the first model to exceed this threshold, with the prior state-of-the-art result at 86.4%. Achieving high performance requires specialist knowledge across many domains (e.g. law, biology, history, etc.), alongside reading comprehension and reasoning. We find Gemini Ultra achieves highest accuracy when used in combination with a chain-of-thought prompting approach ([Wei et al., 2022b](#)) that accounts for model uncertainty. The model produces a chain of thought with k samples, for example 8 or 32. If there is a consensus above a preset threshold (selected based on the validation split), it selects this answer, otherwise it reverts to a greedy sample based on maximum likelihood choice without chain of thought. We refer the reader to appendix for a detailed breakdown of how this approach compares with only chain-of-thought prompting or only greedy sampling.

In mathematics, a field commonly used to benchmark the analytical capabilities of models, Gemini Ultra shows strong performance on both elementary exams and competition-grade problem sets. For the grade-school math benchmark, GSM8K ([Cobbe et al., 2021](#)), we find Gemini Ultra reaches 94.4%

accuracy with chain-of-thought prompting and self-consistency (Wang et al., 2022) compared to the previous best accuracy of 92% with the same prompting technique. Similar positive trends are observed in increased difficulty math problems drawn from middle- and high-school math competitions (MATH benchmark), with the Gemini Ultra model outperforming all competitor models, reaching 53.2% using 4-shot prompting. The model also outperforms the state of the art on even harder tasks derived from American Mathematical Competitions (150 questions from 2022 and 2023). Smaller models perform poorly on this challenging task scoring close to random, but Gemini Ultra can solve 32% of the questions, compared to the 30% solve rate for GPT-4.

Gemini Ultra also excels in coding, a popular use case of current LLMs. We evaluate the model on many conventional and internal benchmarks and also measure its performance as part of more complex reasoning systems such as AlphaCode 2 (see Section 5.1.7 on complex reasoning systems). For example, on HumanEval, a standard code-completion benchmark (Chen et al., 2021) mapping function descriptions to Python implementations, instruction-tuned Gemini Ultra correctly implements 74.4% of problems. On a new held-out evaluation benchmark for python code generation tasks, Natural2Code, where we ensure no web leakage, Gemini Ultra achieves the highest score of 74.9%.

Evaluation on these benchmarks is challenging and may be affected by data contamination. We performed an extensive leaked data analysis after training to ensure the results we report here are as scientifically sound as possible, but still found some minor issues and decided not to report results on e.g. LAMBADA (Paperino et al., 2016). As part of the evaluation process, on a popular benchmark, HellaSwag (Zellers et al., 2019), we find that an additional hundred fine-tuning steps on specific website extracts corresponding to the HellaSwag training set (which were not included in the Gemini model pretraining set) improve the validation accuracy of Gemini Pro to 89.6% and Gemini Ultra to 96.0%, when measured with 1-shot prompting (we measured GPT-4 obtained 92.3% when evaluated 1-shot via the API). This suggests that the benchmark results are susceptible to the pretraining dataset composition. We choose to report HellaSwag decontaminated results only in a 10-shot evaluation setting. We believe there is a need for more robust and nuanced standardized evaluation benchmarks with no leaked data. So, we evaluate Gemini models on several new held-out evaluation datasets that were recently released, such as WMT23 and Math-AMC 2022-2023 problems, or internally generated from non-web sources, such as Natural2Code. We refer the reader to Appendix 10.3 for a comprehensive list of our evaluation benchmarks.

Even so, model performance on these benchmarks gives us an indication of the model capabilities and where they may provide impact on real-world tasks. For example, Gemini Ultra’s impressive reasoning and STEM competencies pave the way for advancements in LLMs within the educational domain<sup>3</sup>. The ability to tackle complex mathematical and scientific concepts opens up exciting possibilities for personalized learning and intelligent tutoring systems.

### 5.1.2. Trends in Capabilities

We investigate the trends in capabilities across the Gemini model family by evaluating them on a holistic harness of more than 50 benchmarks in six different capabilities, noting that some of the most notable benchmarks were discussed in the last section. These capabilities are: “Factuality” covering open/closed-book retrieval and question answering tasks; “Long-Context” covering long-form summarization, retrieval and question answering tasks; “Math/Science” including tasks for mathematical problem solving, theorem proving, and scientific exams; “Reasoning” tasks that require arithmetic, scientific, and commonsense reasoning; “Multilingual” tasks for translation, summarization, and reasoning in multiple languages. Several of these capabilities are targeted by post-training (Section 6). Please see Appendix 10.3 for a detailed list of tasks included for each capability.

---

<sup>3</sup>See demos on website <https://deepmind.google/gemini>.

	Gemini Ultra	Gemini Pro	GPT-4	GPT-3.5	PaLM 2-L	Claude 2	Inflection-2	Grok 1	LLAMA-2
<b>MMLU</b> Multiple-choice questions in 57 subjects (professional & academic) (Hendrycks et al., 2021a)	<b>90.04%</b> CoT@32*	79.13% CoT@8*	87.29% CoT@32 (via API**)	70% 5-shot	78.4% 5-shot	78.5% 5-shot CoT	79.6% 5-shot	73.0% 5-shot	68.0%***
<b>GSM8K</b> Grade-school math (Cobbe et al., 2021)	<b>94.4%</b> Maj1@32	86.5% Maj1@32	86.4% 5-shot (reported)	92.0% SFT & 5-shot CoT	57.1% 5-shot	80.0% 5-shot	88.0% 0-shot	81.4% 8-shot	62.9% 8-shot
<b>MATH</b> Math problems across 5 difficulty levels & 7 subdisciplines (Hendrycks et al., 2021b)	<b>53.2%</b> 4-shot	32.6% 4-shot	52.9% 4-shot (via API**)	34.1% 4-shot (via API**)	34.4% 4-shot	—	34.8%	23.9% 4-shot	13.5% 4-shot
<b>BIG-Bench-Hard</b> Subset of hard BIG-bench tasks written as CoT problems (Srivastava et al., 2022)	<b>83.6%</b> 3-shot	75.0% 3-shot	83.1% 3-shot (via API**)	66.6% 3-shot (via API**)	77.7% 3-shot	—	—	—	51.2% 3-shot
<b>HumanEval</b> Python coding tasks (Chen et al., 2021)	<b>74.4%</b> 0-shot (PT****)	67.7% 0-shot (PT****)	67.0% 0-shot (reported)	48.1% 0-shot	—	70.0% 0-shot	44.5% 0-shot	63.2% 0-shot	29.9% 0-shot
<b>Natural2Code</b> Python code generation. (New held-out set with no leakage on web)	<b>74.9%</b> 0-shot	69.6% 0-shot	73.9% 0-shot (via API**)	62.3% 0-shot (via API**)	—	—	—	—	—
<b>DROP</b> Reading comprehension & arithmetic. (metric: F1-score) (Dua et al., 2019)	<b>82.4</b> Variable shots	74.1 Variable shots	80.9 3-shot (reported)	64.1 3-shot	82.0 Variable shots	—	—	—	—
<b>HellaSwag</b> (validation set) Common-sense multiple choice questions (Zellers et al., 2019)	87.8% 10-shot	84.7% 10-shot	<b>95.3%</b> 10-shot (reported)	85.5% 10-shot	86.8% 10-shot	—	89.0% 10-shot	—	80.0%***
<b>WMT23</b> Machine translation (metric: BLEURT) (Tom et al., 2023)	<b>74.4</b> 1-shot (PT****)	71.7 1-shot	73.8 1-shot (via API**)	—	72.7 1-shot	—	—	—	—

**Table 2 | Gemini performance on text benchmarks with external comparisons and PaLM 2-L.**

\* The model produces a chain of thought with  $k = 8$  or 32 samples, if there is a consensus above a threshold (chosen based on the validation split), it selects this answer, otherwise it reverts to a greedy sample. Further analysis in Appendix 10.2.

\*\* Results self-collected via the API in Nov, 2023.

\*\*\* Results shown use the decontaminated numbers from Touvron et al. (2023b) report as the most relevant comparison to Gemini models which have been decontaminated as well.)

\*\*\*\* PT denotes a post-trained Gemini API model.

We observe consistent quality gains with increased model size in Figure 3, especially in reasoning, math/science, summarization and long-context. Gemini Ultra is the best model across the board for all six capabilities. Gemini Pro, the second-largest model in the Gemini family of models, is also quite competitive while being a lot more efficient to serve.

### 5.1.3. Nano

Bringing AI closer to the user, we discuss the Gemini Nano 1 and Nano 2 models engineered for on-device deployments. These models excel in summarization and reading comprehension tasks with per-task fine-tuning. Figure 3 shows the performance of these pre-trained models in comparison to the much larger Gemini Pro model, while Table 3 dives deeper into specific factuality, coding, Math/Science, and reasoning tasks. Nano-1 and Nano-2 model sizes are only 1.8B and 3.25B parameters respectively. Despite their size, they show exceptionally strong performance on factuality, i.e. retrieval-related tasks, and significant performance on reasoning, STEM, coding, multimodal and

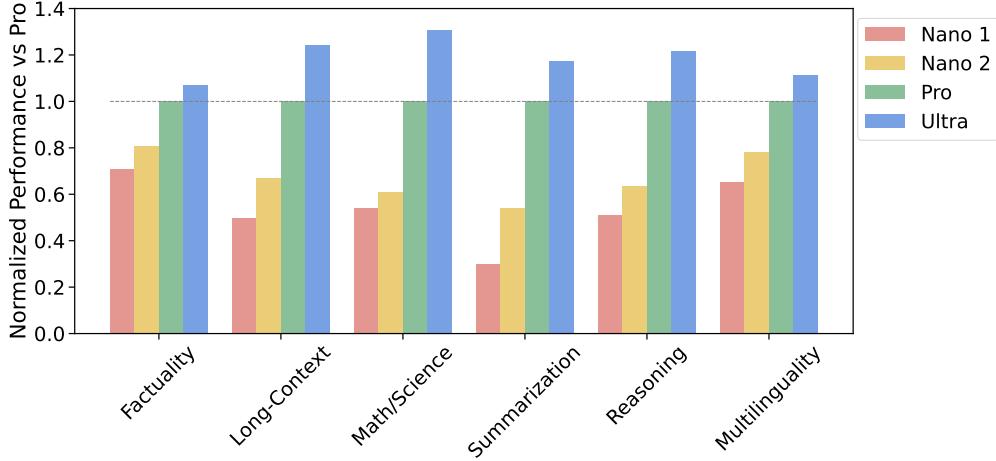


Figure 3 | Language understanding and generation performance of Gemini model family across different capabilities (normalized by the Gemini Pro model).

multilingual tasks. With new capabilities accessible to a broader set of platforms and devices, the Gemini models expand accessibility to everyone.

	Gemini Nano 1		Gemini Nano 2	
	accuracy	normalized by Pro	accuracy	normalized by Pro
BoolQ	71.6	0.81	79.3	0.90
TydiQA (GoldP)	68.9	0.85	74.2	0.91
NaturalQuestions (Retrieved)	38.6	0.69	46.5	0.83
NaturalQuestions (Closed-book)	18.8	0.43	24.8	0.56
BIG-Bench-Hard (3-shot)	34.8	0.47	42.4	0.58
MBPP	20.0	0.33	27.2	0.45
MATH (4-shot)	13.5	0.41	22.8	0.70
MMLU (5-shot)	45.9	0.64	55.8	0.78

Table 3 | Performance of Gemini Nano series on factuality, summarization, reasoning, coding and STEM tasks compared to significantly larger Gemini Pro model.

#### 5.1.4. Multilinguality

The multilingual capabilities of the Gemini models are evaluated using a diverse set of tasks requiring multilingual understanding, cross-lingual generalization, and the generation of text in multiple languages. These tasks include machine translation benchmarks (WMT 23 for high-medium-low resource translation; Flores, NTREX for low and very low resource languages), summarization benchmarks (XLSum, Wikilingua), and translated versions of common benchmarks (MGSM: professionally translated into 11 languages).

##### 5.1.4.1 Machine Translation

Translation is a canonical benchmark in machine learning with a rich history. We evaluated a post-trained Gemini API Ultra model (see Section 6.5.3) on the entire set of language pairs in the WMT 23 translation benchmark in a few-shot setting. Overall, we found that Gemini Ultra (and other Gemini models) performed remarkably well at translating from English to any other language, and surpassed

the LLM-based translation methods when translating out-of-English, on high-resource, mid-resource and low-resource languages. In the WMT 23 out-of-English translation tasks, Gemini Ultra achieved the highest LLM-based translation quality, with an average BLEURT (Sellam et al., 2020) score of 74.8, compared to GPT-4’s score of 73.6, and PaLM 2’s score of 72.2. When averaged across all language pairs and directions for WMT 23, we see a similar trend with Gemini Ultra 74.4, GPT-4 73.8 and PaLM 2-L 72.7 average BLEURT scores on this benchmark.

WMT 23 (Avg BLEURT)	Gemini Ultra	Gemini Pro	Gemini Nano 2	Gemini Nano 1	GPT-4	PaLM 2-L
High Resource	<b>74.2</b>	71.7	67.7	64.1	74.0	72.6
Mid Resource	<b>74.7</b>	71.8	67.0	64.8	73.6	72.7
Out-of-English	<b>74.8</b>	71.5	66.2	65.2	73.6	72.2
Into-English	73.9	72.0	69.0	63.5	<b>74.1</b>	73.4
All languages	<b>74.4</b>	71.7	67.4	64.8	73.8	72.7

Table 4 | Performance of Gemini models on WMT 23 translation benchmark. All numbers with 1-shot.

In addition to the languages and translation tasks above, we also evaluate Gemini Ultra on very low-resource languages. These languages were sampled from the tail of the following language sets: Flores-200 (Tamazight and Kanure), NTREX (North Ndebele), and an internal benchmark (Quechua). For these languages, both from and into English, Gemini Ultra achieved an average chrF score of 27.0 in 1-shot setup, while the next-best model, PaLM 2-L, achieved a score of 25.3.

#### 5.1.4.2 Multilingual Math and Summarization

Beyond translation, we evaluated how well Gemini models perform in challenging tasks across a range of languages. We specifically investigated the math benchmark MGSM (Shi et al., 2023), which is a translated variant of the math benchmark GSM8K (Cobbe et al., 2021). We find Gemini Ultra achieves an accuracy of 79.0%, an advance over PaLM 2-L which scores 74.7%, when averaged across all languages in an 8-shot setup. We also benchmark Gemini models on the multilingual summarization benchmarks – XLSum (Hasan et al., 2021) and WikiLingua (Ladhak et al., 2020). In XLSum, Gemini Ultra reached an average of 17.6 rougeL score compared to 15.4 for PaLM 2. For Wikilingua, Gemini Ultra (5-shot) trails behind PaLM 2 (3-shot) measured in BLEURT score. See Table 5 for the full results. Overall the diverse set of multilingual benchmarks show that Gemini family models have a broad language coverage, enabling them to also reach locales and regions with low-resource languages.

	Gemini Ultra	Gemini Pro	GPT-4	PaLM 2-L
MGSM (8-shot)	<b>79.0</b>	63.5	74.5	74.7
XLsum (3-shot)	<b>17.6</b>	16.2	—	15.4
Wikilingua	48.9	47.8	—	<b>50.4</b>

Table 5 | Performance of Gemini models on multilingual math and summarization.

#### 5.1.5. Long Context

Gemini models are trained with a sequence length of 32,768 tokens and we find that they make use of their context length effectively. We first verify this by running a synthetic retrieval test: we place key-value pairs at the beginning of the context, then add long filler text, and ask for value associated with a particular key. We find that the Ultra model retrieves the correct value with 98% accuracy when queried across the full context length. We further investigate this by plotting the negative log

likelihood (NLL) versus the token index across a held-out set of long documents in Figure 4. We find that the NLL decreases with sequence position up to the full 32K context length. The longer context length of Gemini models enable new use cases such as retrieval over documents and video understanding discussed in Section 5.2.2.

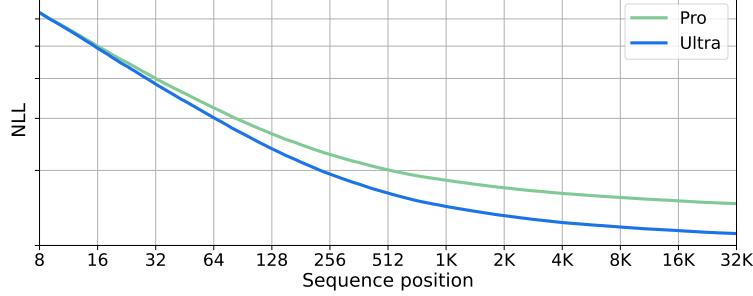


Figure 4 | Negative log likelihood as a function of token index across 32K context length on a held-out set of long documents.

### 5.1.6. Factuality

Factuality ([Maynez et al., 2020](#)) is a key focus of our model’s training and deployment. We evaluate three aspects of factuality for our Gemini API models:

1. **Closed-Book Factuality:** If provided with a fact-seeking prompt without any given source, Gemini API models should not hallucinate incorrect information (see Section 2 of [Roberts et al. \(2020\)](#) for a definition). These prompts can range from information-seeking prompts (e.g. “Who is the prime minister of India?”) to semi-creative prompts that may request factual information (e.g. “Write a 500-word speech in favor of the adoption of renewable energy”).
2. **Attribution:** If instructed to generate a response grounded to a given context, we aim to ensure that Gemini API models produce a response with the highest degree of faithfulness to the context ([Maynez et al., 2020](#); [Rashkin et al., 2023](#)). This may include the summarization of a user-provided source, generating fine-grained citations given a question and provided snippets akin to [Menick et al. \(2022\)](#); [Peng et al. \(2023\)](#), answering questions from a long-form source such as a book ([Mihaylov et al., 2018](#)), and transforming a given source to a desired output (e.g. an email from a portion of a meeting transcript).
3. **Hedging:** If prompted with an input that is “unanswerable”, Gemini API models must acknowledge that it cannot provide a response by hedging to avoid hallucination. These include scenarios where the input prompt contains false-premise questions [see examples in [Hu et al. \(2023\)](#)], the input prompt instructs the model to perform open book QA, but the answer is not derivable from the given context, and so forth.

Factuality is evaluated via human annotators who fact-check each response manually; we report the percentage of factually inaccurate responses as judged by annotators. Attribution is evaluated via human annotators who check for attribution to sources in the prompt for each response manually; the reported metric is AIS ([Rashkin et al., 2023](#)). For hedging, we use an automatic evaluation setup where we measure whether models hedge accurately.

We compare Gemini API Pro with a version without any factuality-focused adaptation in Table 6. We see that the rate of inaccuracy is halved in the factuality set, the accuracy of attribution is increased

by 50% from the attribution set, and the model successfully hedges 70% (up from 0%) in the provided hedging set task.

	Factuality (Inaccurate Rate)	Attribution (AIS)	Hedging (Accuracy)
Gemini API Pro No factuality-focused adaptation	6.7% [5.8%, 7.8%]	40.2% [37.9%, 42.5%]	0%
Gemini API Pro Final stage of post-training	3.8% [3.1%, 4.8%]	60.0% [57.6%, 62.1%]	69.3%

Table 6 | Factuality mitigations: Impact of post-training on the rate of inaccuracy, presence of attribution and the rate of accurate hedging on Gemini API Pro (with corresponding 95% confidence intervals).

### 5.1.7. Complex Reasoning Systems

Gemini models can also be combined with additional techniques such as search and tool-use to create powerful reasoning systems that can tackle more complex multi-step problems. One example of such a system is AlphaCode 2, a new state-of-the-art agent that excels at solving competitive programming problems (Leblond et al, 2023). AlphaCode 2 uses a specialized version of Gemini Pro – tuned on competitive programming data similar to the data used in Li et al. (2022) – to conduct a massive search over the space of possible programs. This is followed by a tailored filtering, clustering and reranking mechanism. Gemini Pro is fine-tuned both to be a coding model to generate proposal solution candidates, and to be a reward model that is leveraged to recognize and extract the most promising code candidates.

AlphaCode 2 is evaluated on Codeforces,<sup>4</sup> the same platform as AlphaCode, on 12 contests from division 1 and 2, for a total of 77 problems. AlphaCode 2 solved 43% of these competition problems, a 1.7x improvement over the prior record-setting AlphaCode system which solved 25%. Mapping this to competition rankings, AlphaCode 2 built on top of Gemini Pro sits at an estimated 85th percentile on average – i.e. it performs better than 85% of entrants. This is a significant advance over AlphaCode, which only outperformed 50% of competitors.

The composition of powerful pre-trained models with search and reasoning mechanisms is an exciting direction towards more general agents; another key ingredient is deep understanding across a range of modalities which we discuss in the next section.

<sup>4</sup><http://codeforces.com/>

## 5.2. Multimodal

Gemini models are natively multimodal. These models exhibit the unique ability to seamlessly combine their capabilities across modalities (e.g. extracting information and spatial layout out of a table, a chart, or a figure) with the strong reasoning capabilities of a language model (e.g. its state-of-art-performance in math and coding) as seen in examples in Figures 5 and 14. The models also show strong performance in discerning fine-grained details in inputs, aggregating context across space and time, and applying these capabilities over a temporally-related sequence of video frames and/or audio inputs.

The sections below provide more detailed evaluation of the model across different modalities (image, video, and audio), together with qualitative examples of the model’s capabilities for image generation and the ability to combine information across different modalities.

### 5.2.1. Image Understanding

We evaluate post-trained Gemini API models on four different capabilities: high-level object recognition using captioning or question-answering tasks such as VQAv2; fine-grained transcription using tasks such as TextVQA and DocVQA requiring the model to recognize low-level details; chart understanding requiring spatial understanding of input layout using ChartQA and InfographicVQA tasks; and multimodal reasoning using tasks such as Ai2D, MathVista and MMMU. For zero-shot QA evaluation, the model is instructed to provide short answers aligned with the specific benchmark. All numbers are obtained using greedy sampling and without any use of external OCR tools.

	Gemini Ultra (pixel only)	Gemini Pro (pixel only)	Gemini Nano 2 (pixel only)	Gemini Nano 1 (pixel only)	GPT-4V	Prior SOTA
<b>MMMU (val)</b> Multi-discipline college-level problems (Yue et al., 2023)	<b>59.4%</b> pass@1  <b>62.4%</b> Maj1@32	47.9%	32.6%	26.3%	56.8%	56.8% GPT-4V, 0-shot
<b>TextVQA (val)</b> Text reading on natural images (Singh et al., 2019)	<b>82.3%</b>	74.6%	65.9%	62.5%	78.0%	79.5% Google PaLI-3, fine-tuned
<b>DocVQA (test)</b> Document understanding (Mathew et al., 2021)	<b>90.9%</b>	88.1%	74.3%	72.2%	88.4% (pixel only)	88.4% GPT-4V, 0-shot
<b>ChartQA (test)</b> Chart understanding (Masry et al., 2022)	<b>80.8%</b>	74.1%	51.9%	53.6%	78.5% (4-shot CoT)	79.3% Google DePlot, 1-shot PoT (Liu et al., 2023)
<b>InfographicVQA (test)</b> Infographic understanding (Mathew et al., 2022)	<b>80.3%</b>	75.2%	54.5%	51.1%	75.1% (pixel only)	75.1% GPT-4V, 0-shot
<b>MathVista (testmini)</b> Mathematical reasoning (Lu et al., 2023)	<b>53.0%</b>	45.2%	30.6%	27.3%	49.9%	49.9% GPT-4V, 0-shot
<b>Ai2D (test)</b> Science diagrams (Kembhavi et al., 2016)	<b>79.5%</b>	73.9%	51.0%	37.9%	78.2%	81.4% Google PaLI-X, fine-tuned
<b>VQAv2 (test-dev)</b> Natural image understanding (Goyal et al., 2017)	<b>77.8%</b>	71.2%	67.5%	62.7%	77.2%	86.1% Google PaLI-X, fine-tuned

Table 7 | **Image understanding** Gemini Ultra consistently outperforms existing approaches even in zero-shot, especially for OCR-related image understanding tasks for natural images, text, documents, and figures without using any external OCR engine (‘pixel only’). Many existing approaches fine-tune on the respective tasks, highlighted in gray, which makes the comparison with 0-shot not apples-to-apples.

We find that Gemini Ultra is state of the art across a wide range of image-understanding benchmarks in Table 7. It achieves strong performance across a diverse set of tasks such as answering questions on natural images and scanned documents as well as understanding infographics, charts and science diagrams. When compared against publicly reported results from other models (most notably GPT-4V), the Gemini model is better in zero-shot evaluation by a significant margin. It also exceeds several existing models that are specifically fine-tuned on the benchmark’s training sets for the majority of tasks. The capabilities of the Gemini models lead to significant improvements in the state of the art on academic benchmarks like MathVista (+3.1%)<sup>5</sup> or InfographicVQA (+5.2%).

MMMU (Yue et al., 2023) is a recently released evaluation benchmark, which consists of questions about images across 6 disciplines with multiple subjects within each discipline that require college-level knowledge to solve these questions. Gemini Ultra achieves the best score on this benchmark advancing the state-of-the-art result by more than 5 percentage points and outperforms the previous best result in 5 of 6 disciplines (see Table 8), thus showcasing its multimodal reasoning capabilities.

MMMU (val)	Gemini Ultra (0-shot)	GPT-4V (0-shot)	
	Maj@32	pass@1	pass@1
Art & Design	<b>74.2</b>	70.0	65.8
Business	<b>62.7</b>	56.7	59.3
Science	49.3	48.0	<b>54.7</b>
Health & Medicine	<b>71.3</b>	67.3	64.7
Humanities & Social Science	<b>78.3</b>	78.3	72.5
Technology & Engineering	<b>53.0</b>	47.1	36.7
Overall	<b>62.4</b>	59.4	56.8

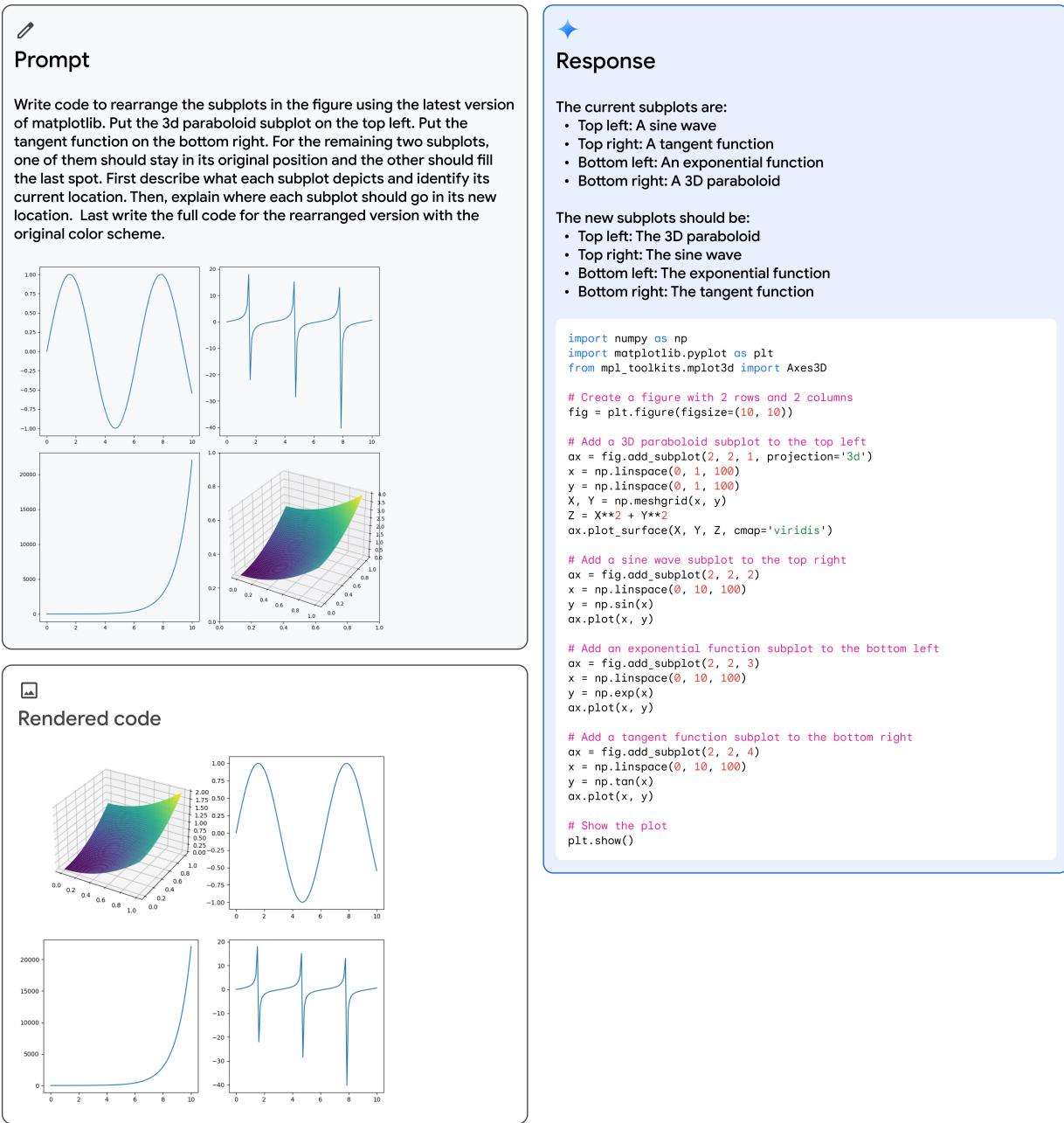
Table 8 | **Gemini Ultra performance on the MMMU benchmark** (Yue et al., 2023) per discipline. Each discipline covers multiple subjects, requiring college-level knowledge and complex reasoning.

Gemini models are also capable of operating across modalities and a diverse set of global languages simultaneously, both for image understanding tasks (e.g., images containing text in Icelandic) and for generation tasks (e.g., generating image descriptions for a wide range of languages). We evaluate the performance of generating image descriptions on a selected subset of languages in the Crossmodal-3600 (XM-3600) benchmark in a 4-shot setting, using the Flamingo evaluation protocol (Alayrac et al., 2022), without any fine-tuning for all models. As shown in Table 9, Gemini models achieve a significant improvement over the existing best model, Google PaLI-X.

XM-3600 (CIDEr)	Gemini Ultra 4-shot	Gemini Pro 4-shot	Google PaLI-X 4-shot
English	86.4	<b>87.1</b>	77.8
French	<b>77.9</b>	76.7	62.5
Hindi	<b>31.1</b>	29.8	22.2
Modern Hebrew	<b>54.5</b>	52.6	38.7
Romanian	<b>39.0</b>	37.7	30.2
Thai	<b>86.7</b>	77.0	56.0
Chinese	<b>33.3</b>	30.2	27.7
Average (of 7)	<b>58.4</b>	55.9	45.0

Table 9 | **Multilingual image understanding** Gemini models outperform existing models in captioning images in many languages when benchmarked on a subset of languages in XM-3600 dataset (Thapliyal et al., 2022).

<sup>5</sup>MathVista is a comprehensive mathematical reasoning benchmark consisting of 28 previously published multimodal datasets and three newly created datasets. Our MathVista results were obtained by running the [MathVista authors’ evaluation script](#).



**Figure 5 | Using Gemini models’ multimodal reasoning capabilities** to generate matplotlib code for rearranging the subplots. The multimodal prompt is shown at the top-left in gray. Gemini Ultra’s response, including its generated code, is shown in the right column in blue. The bottom left figure shows rendered version of the generated code. Successfully solving this task shows the model’s capability to combine several capabilities: (1) **recognition** of the functions depicted in the plots; (2) **inverse graphics** to infer the code that would have generated the subplots; (3) **instruction-following** to put subplots in their desired positions; and (4) **abstract reasoning** to infer that the exponential plot must stay in its original place, because the sine plot must move out of the way for the 3-dimensional plot.

Qualitative evaluation in Figure 5 illustrates an example of Gemini Ultra’s multimodal reasoning capabilities. The model is required to solve the task of generating matplotlib code that would rearrange

a set of subplots provided by the user. The model output shows that it successfully solves this task combining multiple capabilities of understanding the user plot, inferring the code required to generate it, following user instructions to put subplots in their desired positions, and abstract reasoning about the output plot. This highlights Gemini Ultra’s native multimodality and alludes to its more complex reasoning abilities across interleaved sequences of image and text. We refer the reader to the appendix for more qualitative examples.

### 5.2.2. Video Understanding

Understanding video input is an important step towards a useful generalist agent. We measure the video understanding capability across several established benchmarks that are held-out from training. These tasks measure whether the model is able to understand and reason over a temporally-related sequence of frames. For each video task, we sample 16 equally-spaced frames from each video clip and feed them to the Gemini models. For the YouTube video datasets (all datasets except NextQA and the Perception test), we evaluate the Gemini models on videos that were still publicly available in the month of November, 2023.

Gemini Ultra achieves state-of-the-art performance on various few-shot video captioning tasks as well as zero-shot video question answering tasks as shown in Table 10. This demonstrates its capability of strong temporal reasoning across several frames. Figure 23 in the appendix provides a qualitative example of understanding the video of the ball-striking mechanics of a soccer player and reasoning about the player can improve their game.

Task	Gemini Ultra	Gemini Pro	Few-shot SoTA
<b>VATEX (test)</b> English video captioning (Wang et al., 2019)	<b>62.7</b> 4-shots	57.4 4-shots	56.0 DeepMind Flamingo, 4-shots
<b>VATEX ZH (test)</b> Chinese video captioning (Wang et al., 2019)	<b>51.3</b> 4-shots	50.0 4-shots	—
<b>YouCook2 (val)</b> English cooking video captioning (Zhou et al., 2018)	<b>135.4</b> 4-shots	123.2 4-shots	74.5 DeepMind Flamingo, 4-shots
<b>NextQA (test)</b> Video question answering (Xiao et al., 2021)	<b>29.9</b> 0-shot	28.0 0-shot	26.7 DeepMind Flamingo, 0-shot
<b>ActivityNet-QA (test)</b> Video question answering (Yu et al., 2019)	<b>52.2</b> 0-shot	49.8 0-shot	45.3 Video-LLAVA, 0-shot
<b>Perception Test MCQA (test)</b> Video question answering (Pătrăucean et al., 2023)	<b>54.7</b> 0-shot	51.1 0-shot	46.3 SeViLA (Yu et al., 2023), 0-shot

Table 10 | **Few-shot video understanding across tasks and languages** on selected academic benchmarks. The reported metric is CIDEr for video captioning, WUPS for NextQA, and top-1 accuracy for the Perception Test and ActivityNet-QA. For ActivityNet-QA, we use the Video-LLAVA (Lin et al., 2023) evaluation protocol.

### 5.2.3. Image Generation

Gemini models are able to output images natively, without having to rely on an intermediate natural language description that can bottleneck the model’s ability to express images. This uniquely enables the model to generate images with prompts using interleaved sequences of image and text in a

few-shot setting. For example, the user might prompt the model to design suggestions of images and text for a blog post or a website (see Figure 12 in the appendix).

Figure 6 shows an example of image generation in 1-shot setting. Gemini Ultra model is prompted with one example of interleaved image and text where the user provides two colors (blue and yellow) and image suggestions of creating a cute blue cat or a blue dog with yellow ear from yarn. The model is then given two new colors (pink and green) and asked for two ideas about what to create using these colors. The model successfully generates an interleaved sequence of images and text with suggestions to create a cute green avocado with pink seed or a green bunny with pink ears from yarn.

 **Prompt**

Give me two ideas that I could do with these 2 colors

Colors: I see blue and yellow yarn

How about a cute blue cat?

Or a blue dog that would also have a yellow ear?

Give me two ideas that I could do with these 2 colors






 **Response**

Colors: I see green and pink yarn

Idea 1: How about a green avocado with pink seed?

Idea 2: Or a green bunny with pink ears?




Figure 6 | **Image Generation.** Gemini models can output multiple images interleaved with text given a prompt composed of image and text. In the left figure, Gemini Ultra is prompted in a 1-shot setting with a user example of generating suggestions of creating cat and dog from yarn when given two colors, blue and yellow. Then, the model is prompted to generate creative suggestions with two new colors, pink and green, and it generates images of creative suggestions to make a cute green avocado with pink seed or a green bunny with pink ears from yarn as shown in the right figure.

### 5.2.4. Audio Understanding

We evaluate the Gemini Nano-1 and Gemini Pro models on a variety of public benchmarks and compare it with Universal Speech Model (USM) (Zhang et al., 2023) and Whisper (large-v2 (Radford et al., 2023) or large-v3 (OpenAI, 2023) as indicated). These benchmarks include automatic speech recognition (ASR) tasks such as FLEURS (Conneau et al., 2023), VoxPopuli, (Wang et al., 2021), Multi-lingual Librispeech (Pratap et al., 2020), as well as the speech translation task CoVoST 2, translating different languages into English (Wang et al., 2020). We also report on an internal benchmark YouTube test set. ASR tasks report a word error rate (WER) metric, where a lower number is better. Translation tasks report a BiLingual Evaluation Understudy (BLEU) score, where a higher number is better. FLEURS is reported on 62 languages that have language overlap with the training data. Four segmented languages (Mandarin, Japanese, Korean and Thai) report character error rate (CER), instead of WER, similar to Whisper (Radford et al., 2023).

Table 11 indicates that our Gemini Pro model significantly outperforms the USM and Whisper models across all ASR and AST tasks, both for English and multilingual test sets. Note that there is a large gain in FLEURS, compared to USM and Whisper, as our model is also trained with the FLEURS training dataset. However, training the same model without FLEURS dataset results in a WER of 15.8, which still outperforms Whisper. Gemini Nano-1 model also outperforms both USM and Whisper on all datasets except FLEURS. Note that we did not evaluate Gemini Ultra on audio yet, though we expect better performance from increased model scale.

	Task	Metric	Gemini Pro	Gemini Nano-1	Whisper (OpenAI, 2023; Radford et al., 2023)	USM (Zhang et al., 2023)
Automatic Speech Recognition	<b>YouTube</b> (en-us)	WER (↓)	<b>4.9%</b>	5.5%	6.5% (v3)	6.2%
	<b>Multilingual Librispeech</b> (en-us) (Pratap et al., 2020)	WER (↓)	<b>4.8%</b>	5.9%	6.2% (v2)	7.0 %
	<b>FLEURS</b> (62 lang) (Conneau et al., 2023)	WER (↓)	<b>7.6%</b>	14.2%	17.6% (v3)	11.8%
	<b>VoxPopuli</b> (14 lang) (Wang et al., 2021)	WER (↓)	<b>9.1%</b>	9.5%	15.9% (v2)	13.4%
	<b>CoVoST 2</b> (21 lang) (Wang et al., 2020)	BLEU (↑)	<b>40.1</b>	35.4	29.1 (v2)	30.7

Table 11 | Speech evaluation results on selected benchmarks for ASR and AST. For ASR, the reported metric is WER where lower is better. For AST, the reported metric is BLEU where higher is better.

Table 12 shows further error analysis with USM and Gemini Pro. We find that Gemini Pro produces more understandable responses, particularly on rare words and proper nouns.

Domain	Truth	USM	Gemini Pro	Wav
Fleurs	Scotturb bus 403 travels regularly to Sintra, stopping at Cabo da Roca.	Scotboard bus four3 traversed regularly to Centra stopping at Cabo de Roga.	Scotturb bus 403 travels regularly to Sintra, stopping at Cabo da Roca.	
Fleurs	The archipelago lies 120 km north of the Peninsula. The largest is King George Island, with the settlement of Villa Las Estrellas.	The archipelago lines 120 km north of peninsula. The largest is Kingurch island with the settlement of Cua Losas.	The archipelago lies 120 km north of the Peninsula. The largest is King George Island, with the settlement of Villa Las Estrellas.	

Table 12 | Qualitative examples for the ASR task in the benchmark. Incorrect transcriptions are highlighted in red.

### 5.2.5. Modality Combination

Multimodal demonstrations often include a combination of text interleaved with a single modality, usually images. We demonstrate the ability to process a sequence of audio and images natively.

Consider a cooking scenario about making an omelet where we prompt the model with a sequence of audio and images. Table 13 indicates a turn-by-turn interaction with the model, providing pictures and verbally asking questions about the next steps for cooking an omelet. We note that the model response text is reasonably accurate, and shows that model processes fine-grained image details to evaluate when the omelet is fully cooked. See demo on the [website](#).

Input Image	Input Audio (transcribed)	Model Response: Text
	🔊 What's the first step to make a veggie omelet with these ingredients?	Crack the eggs into a bowl and whisk them.
	🔊 Thank you for the instructions. I started making my omelet, does it look ready now?	It looks like it's almost ready. You can flip it over to cook the other side.
(No image - it's a follow up on the previous question)	🔊 Why is it not ready?	It's not ready because the eggs are still runny.
	🔊 What about now?	It looks ready now. You can take it off the heat and serve it.

Table 13 | Audio-visual qualitative example showcasing the ability of Gemini models to process interleaved sequences of text, vision, and audio, as well as reason across modalities. This example inputs interleaved images and audio from the user in a cooking scenario. The user prompts the model for instructions to make an omelet and to inspect whether it is fully cooked.

## 6. Post-Training Models

After large-scale pre-training, we apply *post-training*, where one trains on top of a pre-trained model in order to extend the model’s proficiency and to enable a wide variety of capabilities. Namely, we seek to improve overall quality, enhance target capabilities such as coding and multilingual, and ensure alignment and safety criteria are met. We discuss our approach to post-training in this section, highlighting common and distinct aspects of the Gemini Apps and Gemini API model variants.

### 6.1. Gemini Apps: Gemini and Gemini Advanced

Gemini and Gemini Advanced offer direct access to Google’s family of AI models, consisting of the core post-trained Gemini Apps models and the system around it. These models are created by applying specialized post-training on top of Gemini pre-trained models: currently, Gemini gives access to Pro 1.0 and Gemini Advanced gives access to Ultra 1.0. Beyond the core models, the system determines how the models interact with external tools (such as Google Flights, Maps, and Google Workspace), and how to generate responses (filtering, ranking, and streaming). As an area, conversational AI presents several challenges, including: How to understand users’ requests across multi-turn interactions? How to make sure responses are safe, factually grounded, and helpful? How to help users accomplish tasks by using tools external to the models? We discuss how we approach these challenges in the following sections.

### 6.2. Gemini APIs: Google AI Studio and Cloud Vertex AI

Our developer-focused Gemini API models are designed to support both conversational and non-conversational use cases. These models are available through Google AI Studio and Cloud Vertex AI through an easy to use API. Google AI Studio is a free, web-based developer tool to prototype and launch apps quickly with an API key. Vertex AI is a comprehensive AI platform that enables developers to leverage Gemini API models with varied tooling, fully-managed infrastructure, and built-in enterprise security and privacy settings. Gemini APIs make it easy to integrate Gemini API models into any production product or workflow, empowering developers to build applications that can reason across different modalities.

### 6.3. Post-Training Methods & Data

Post-training Gemini models to produce Gemini API and Apps variants involves several stages; see Figure 7. Careful data curation is critical for all stages. First, we collect a diverse set of prompts that are representative of real-world use cases. Second, we apply supervised fine-tuning (SFT) on demonstration data of what the model’s output should be for a given prompt ([Mishra et al., 2021](#); [Ouyang et al., 2022](#); [Wei et al., 2022a](#)). Third, we further collect different possible responses to a given prompt, and collect feedback data over these to train a Reward Model (RM). Finally, using the trained RM, a Reinforcement Learning from Human Feedback (RLHF) stage ([Bai et al., 2022a](#)) is applied to further align the model’s outputs with human preferences. We discuss our methods in more detail below:

**(1) Prompt Data Collection:** A prompt is a user’s input to the model. As well as the most recent user input, this can also include previous user-model interactions. We curate datasets of target prompts. The datasets serve as the basis for our demonstration and feedback data collections, and they are used directly during reinforcement learning. It is important to cover a diverse set of crucial use cases and in both single-turn and multi-turn formats. Data sources include vendor-created data, third-party licensed sources, and synthetic approaches.

**(2) SFT on Demonstration Data:** SFT trains the model to output a desired target response given a prompt. Our Demonstration Data target responses can be directly written by a human expert, or generated by a model and in some cases revised or reviewed by a human. Additionally, we use data analysis tools and heuristics to ensure high data diversity across capabilities, use cases, and semantic clusters.

**(3) RM Training on Feedback Data:** We further collect Feedback Data, for which human raters provide feedback such as relative preferences over candidate responses and feedback regarding individual responses to a given prompt. For many capabilities, rating relative preferences is an easier task than demonstrating an ideal response. Feedback data are collected across creativity, safety, factuality, other capabilities, and other target criteria. We found that the utility of the resulting human feedback data greatly depends on the prompt selection and the sampling strategy used to produce candidate responses. We use this data to train RMs to output rewards that align with human preferences as closely as possible.

**(4) RLHF:** Applying reinforcement learning from human feedback (RLHF) to our models provides further gains over SFT alone. Our approach creates an iterative process in which RL continually pushes the boundaries of the RM, while the RM is continuously improved through evaluation and data collection, leading to progressive improvements in both.

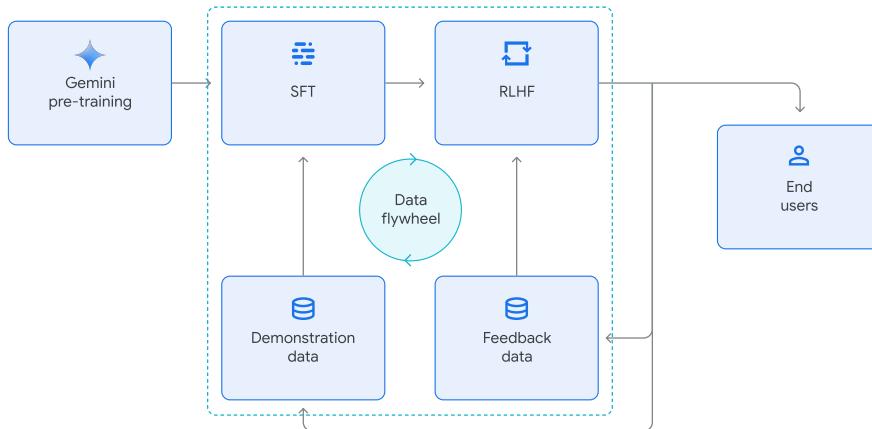


Figure 7 | **Modeling overview.** Post-training utilizes an optimized data flywheel in order to acquire human-AI feedback and continually improve on key areas. The data mixtures for supervised fine-tuning, reward modeling, and reinforcement learning serve as the foundation for our models.

#### 6.4. Evaluation

Evaluation of human preferences over model outputs provides critical signals for measuring performance. As part of our development process, we conduct human evaluation extensively across targeted capabilities. Human evaluation is instantiated as side-by-side blind evaluations where human raters judge responses of two models to the same prompt, as single-response ratings for certain capabilities, and as online testing. In addition, we build models for automated evaluation that faithfully imitate human preferences in order to guide development and continuously monitor online performance.

#### 6.5. Model Capabilities

Beyond the general post-training outlined above, we apply techniques to improve a set of key capabilities. These capabilities cover a range of use cases inspired by current user needs and research-inspired

future applications. We outline capability examples not detailed in previous sections below. The post-training recipes are carefully designed to balance multiple objectives, including creativity, factuality, safety and more (Bai et al., 2022b; Thoppilan et al., 2022). We have a particular focus on safety and alignment, and hence address this in a further dedicated section.

### 6.5.1. Instruction Following

Following a user’s prompt accurately is a fundamental capability for LLMs, especially as these models become more sophisticated and are presented with increasingly complex user prompts. User prompts vary in granularity, specificity, and requirements (e.g., content, format, length). Individual instructions can also be ambiguous, optional, or even impossible or undesirable to satisfy (He et al., 2023; Xu et al., 2023).

We improve Gemini Apps and Gemini API models’ instruction following (IF) abilities by collecting data for a diverse set of instruction following categories. For instructions that are verifiable programmatically such as word count, we generate synthetic data via prompting and response editing to ensure that such instructions are satisfied.

**Complex prompts evaluation:** We investigate performance on complex prompts containing multiple instructions using a fine-grained evaluation method that assesses how well models adhere to each instruction. Human raters are presented with a prompt-response pair and a list of the individual (sub)-instructions contained in the prompt. Each prompt may have anywhere from one to dozens of individual instructions, and the annotators are tasked with determining whether each instruction is followed (or not) by the response.

Table 14 reports results on an internal dataset of prompts with instructions of varying complexity that encompass a wide range of instructions and are designed to be challenging for LLMs. We report two metrics: per-instruction accuracy (the percentage of sub instructions in the eval set that are followed), and full-response accuracy (the percentage of eval set prompts where all sub-instructions are followed).

	Post-trained PaLM 2	Gemini (with Pro)	Gemini Advanced (with Ultra)
Per-instruction accuracy	59.5±3.0%	77.8±2.0%	87.4±1.4%
Full-response accuracy	25.5±3.3%	38.5±3.6%	54.1±3.7%

Table 14 | Performance of Gemini on our complex prompts instruction-following internal benchmark.

Gemini Advanced (with Ultra) achieves an average per-instruction accuracy close to 90%, representing a significant improvement over Gemini (with Pro) and a post-trained PaLM 2 model. We find that the sub-instructions that aren’t followed are well-distributed across responses. As a result Gemini Advanced’s full-response accuracy is lower, at around 54%. This indicates that there is further headroom for models to fully satisfy all instructions.

### 6.5.2. Tool Use

By training LLMs to use tools, we greatly expand LLM capabilities beyond their internal knowledge. We treat tool use for both Gemini Apps and Gemini API models as a code generation problem, leveraging the base model’s preexisting strong coding capabilities. Every tool invocation is represented as a code block in which tool calls are invoked. This process allows the model to both compose multiple tools in each code block, as well as observe and react to the results of tool execution. At inference time, to generate a response to a user prompt, our system executes the loop shown in Figure 8, where sampling from the LLM and execution of tool code work together to create a final response.

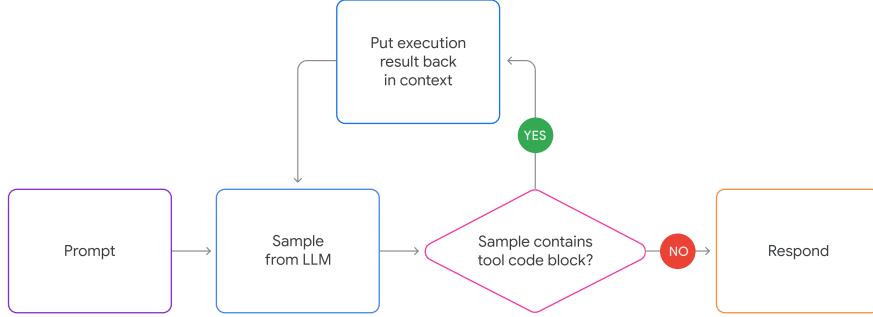


Figure 8 | A Gemini tool-use control loop.

**Gemini Apps models:** Gemini draws on a range of tools via Gemini Extensions, including Google Workspace, Google Maps, YouTube, Google Flights, and Google Hotels. These tool-use capabilities also enable Gemini to be integrated as part of Gmail, Docs, Slides, Sheets and more. We are aiming to bring further tool-use capabilities in order to both enhance Gemini models and integrate Gemini models into further products.

We created an internal benchmark to assess Gemini performance on tasks that may benefit from access to these extensions. This benchmark measures human preference in domains such as travel planning and video discovery. We find models equipped with tools are preferred on this set 78% of the time over models without tools (excluding ties).

**Gemini API models:** We have found that fine-tuning Gemini API models is very effective at teaching the model tool-use behaviors. Furthermore, training models to use programming and search as tools leads to improved performance on a range of academic benchmarks. In Table 15, we compare tool-use models fine-tuned from an early version of Gemini API Pro against equivalent models that do not use tools.

	Mathematical Reasoning		Factuality & Knowledge Retrieval	
	GSM8K Cobbe et al. (2021)	MATH Hendrycks et al. (2021b)	NQ Kwiatkowski et al. (2019b)	Realtime QA Kasai et al. (2022a)
Gemini API Pro with tools	80.1%	41.8%	68.0%	70.8%
Gemini API Pro without tools	69.7%	30.7%	59.0%	39.2%

Table 15 | Comparison between Gemini API tool-use models and comparable models that do not use tools. Gemini API Pro without tools is an early version of our Pro model trained without tool-use data. Gemini API Pro with tools is the same model fine-tuned with tool-use data.

### 6.5.3. Multilinguality

Multilinguality is critical to make sure Gemini models effectively support a wide range of languages. We discuss our key approaches for Gemini Apps and Gemini API models respectively below.

**Gemini Apps models:** Scaling Gemini from English to 40+ languages imposed research challenges in data quality. We leverage abundant high-quality English data by localization to native cultures (e.g., “president of the United States” -> “日本の首相”).

Table 16 shows the performance of Gemini (with Pro) on 5 languages compared to Bard with

an older post-training recipe and based on PaLM 2. For side-by-side comparisons between a model A and a model B, we calculate a metric called SxS score. Each rating is converted to an ordinal value centered at 0: ratings preferring A are positive and ratings preferring B are negative over a scale between -1.5 and 1.5. The converted values are averaged to return the SxS score. Intuitively, a positive SxS score indicates the extent to which model A is preferred over model B. Here, we find quality improved by more than 0.1 SxS score for all five languages. Coding and reasoning gains from Gemini Pro are preserved across languages.

Language	Quality SxS	Coding MBPP Pass@1 <a href="#">Austin et al. (2021)</a>	Reasoning MMLU <a href="#">Hendrycks et al. (2021a)</a>
ja-JP	+0.14	+22.2%	+3.6%
pt-BR	+0.17	+23.2%	+5.2%
de-DE	+0.1	+21.4%	+7.5%
es-419	+0.12	+22.8%	+9.3%
it-IT	+0.13	+13.8%	+7.5%

Table 16 | Multilingual performance of Gemini (with Pro) compared to Gemini with an older post-training recipe and PaLM 2.

**Gemini API models:** Similar to Gemini Apps models, we train Gemini API models on additional multilingual post-training data, effectively adapting the original English model for use in various languages. We experiment with both human-generated non-English prompt-response pairs as well as automatically translated pairs. For the latter, we leverage abundant high-quality English demonstration data by translation. We ensure the quality of such translated data by translationability filtering and response rating by humans.

*Translatability Filtering:* Not all prompt-response pairs make sense when automatically translated, and may require expensive localization instead. Example prompts of this type (responses omitted for space) include:

- (strict word requirements) Write a 1000 word essay about world peace.
- (too English centric) Write a poem in iambic pentameter about apples.
- (too Latin-script centric) What is a word with 1 E, 2 As, and 1 U?

*Translation Quality Validation:* Each translated prompt-response pair was rated for translation quality by at least 3 human raters, and was kept in the final mixture if the majority of raters rated it as accurate. Section 5.1.4 reports evaluations of the multilingual capabilities of post-trained Gemini API models.

#### 6.5.4. Multimodal Vision

Multimodal post-training enhances the capabilities of our natively multimodal Gemini models for a wide range of useful applications. In the following, we discuss how image understanding ability is incorporated into Gemini Apps and Gemini API models. For this evaluation, we further train both of these Gemini model variants on a mixture of text data and expert curated image-text data over several vertically-defined multimodal use cases

**Gemini Apps models:** We empower Gemini and Gemini Advanced with image understanding capabilities by fine-tuning pre-trained Gemini models on a mixture of text-only and image-text data. Careful balancing of text and multimodal data ensures the model develops robust image understanding without adversely affecting the quality of the text-only interactions. To assess our

models, we compile a dataset of human-curated and synthetic image-text prompts and responses, spanning various categories and difficulty levels. This dataset facilitates human evaluation for model comparison and selection.

We find that introducing this image-text data preserves Gemini Apps model quality on text-only tasks, with a SxS score on text-only tasks of  $+0.01 \pm 0.01$  for a Gemini Apps Pro model trained on this data versus an equivalent model trained only on text data. In addition, post-training via RLHF improves performance on multimodal tasks, with a SxS score on image-understanding tasks of  $+0.223 \pm 0.06$  for a Gemini Apps Pro model post-trained with SFT & RLHF vs SFT alone.

**Gemini API models:** We evaluate the impact of post-training via SFT on Gemini API models' multimodal vision performance by tracking the performance of both pre-trained models and post-trained Gemini API Vision models on a series of standard benchmarks. These post-trained results have already been given in Table 7, in Table 17 we further report the difference in performance between pre-trained and post-trained Gemini API models.

	Gemini Ultra Pre-trained only 0-shot (pixel only)	Gemini API Ultra 0-shot (pixel only)	Gemini Ultra pre- to post-trained improvement
<b>MMMU (val)</b> Multi-discipline college-level problems (Yue et al., 2023)	n/a	59.4% pass@1  62.4% Maj1@32	n/a
<b>TextVQA (val)</b> Text reading on natural images (Singh et al., 2019)	81.4%	82.3%	+0.9%
<b>DocVQA (test)</b> Document understanding (Mathew et al., 2021)	90.1%	90.9%	+0.8%
<b>ChartQA (test)</b> Chart understanding (Masry et al., 2022)	80.8%	80.8%	0.0%
<b>InfographicVQA (test)</b> Infographic understanding (Mathew et al., 2022)	77.9%	80.3%	+2.4%
<b>MathVista (testmini)</b> Mathematical reasoning (Lu et al., 2023)	n/a	53.0%	n/a
<b>AI2D (test)</b> Science diagrams (Kembhavi et al., 2016)	76.6%	79.5%	+2.9%
<b>VQAv2 (test-dev)</b> Natural image understanding (Goyal et al., 2017)	74.5%	77.8%	+3.3%

Table 17 | **Post-trained model image understanding** Post-training improves image understanding capabilities of Gemini API Ultra over the base pre-trained model. Comparisons of Gemini API Ultra to other models on these benchmarks are given in Table 7.

The results indicate that the pre-trained model already has high performance across the capabilities represented by these benchmarks, in line with previous observations. However, the post-training SFT stage used for the Gemini API Vision models succeeds in improving the performance over several of these benchmarks (InfographicVQA, AI2D, VQAv2), most likely due to the model's increased instruction-following capabilities that succeed in aligning the model output style with that of the golden references.

### 6.5.5. Coding

Despite the strong coding benchmark performance of the base model, post-training data still provides a significant boost to both code quality and code correctness. This highlights the benefit of high-quality demonstration data and feedback data for coding use cases. Gemini Apps and Gemini API models use a combination of human and synthetic approaches to collect such data.

We evaluate our Gemini Apps models' coding performance on a set of internally curated prompts, distributed across code use cases and languages. Table 18 reports SxS scores, where Gemini (with Pro) significantly improves upon Bard with an older post-training recipe and based on PaLM 2. Gemini Advanced (with Ultra) further improves upon Gemini (with Pro).

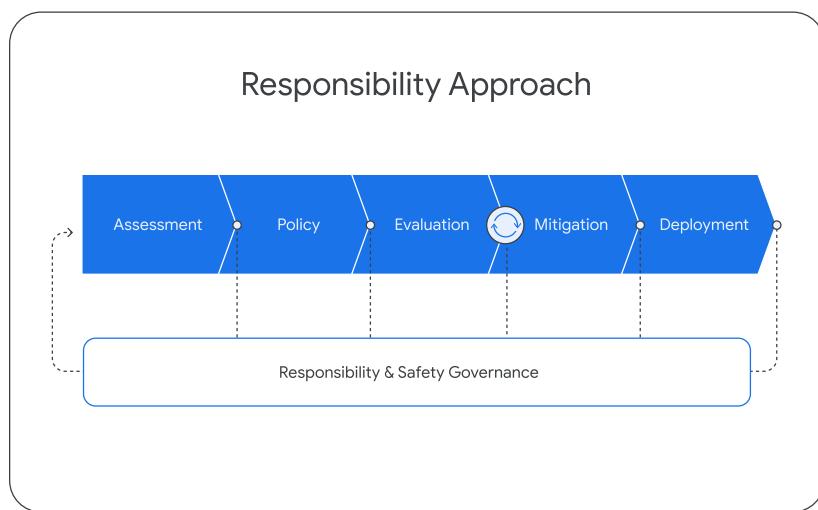
Side A	Side B	SxS score
Gemini (with Pro)	Bard (PaLM 2, Sept. 2023)	0.19±0.03
Gemini Advanced (with Ultra)	Gemini (with Pro)	0.13± 0.02

Table 18 | SxS comparisons of Gemini models on an internal coding benchmark.

For the coding capabilities of post-trained Gemini API Models, see Table 2 which reports their academic benchmark performance.

## 7. Responsible Deployment

During the development of Gemini models, we follow a structured approach to responsible deployment to identify, measure, and manage foreseeable downstream societal impacts of our models, in line with previous releases of Google's AI technology ([Kavukcuoglu et al., 2022](#)). Throughout the lifecycle of a project, we follow the structure below. This section provides more detail about our approach and includes key findings where available. We are committed to ongoing transparency and will continue to provide updated information on our approach and testing in upcoming reports.



### 7.1. Impact Assessment

At Google we apply an impact assessment framework throughout the product development lifecycle related to Google's AI Principles ([Google, 2023](#)). This means we assess the risk and impact of AI models we're building at both a model-level (e.g. for Gemini API Ultra 1.0, as deployed on Cloud

Studio or Vertex AI), and once embedded within a broader product or service (e.g. for Gemini Advanced).

### 7.1.1. Model Assessment

We conduct model impact assessments to identify, assess, and document societal benefits and harms associated with the capabilities of Gemini models. Our impact assessments for Gemini API models describe downstream benefits and risks that we identify, spanning across the models' modalities (text-to-text; image-to-text; and video-to-text). Model impact assessments are conducted by the Google DeepMind Responsible Development and Innovation team, and are reviewed by the Google DeepMind Responsibility and Safety Council. We draw from various sources in producing impact assessments, including a wide range of literature, external expertise, and our in-house ethics and safety research.

Gemini models introduce various benefits to people and society. Gemini models' various modalities, including language, image and video understanding, can help users process information more efficiently, for example through content summarisation. These efficiency benefits can apply to commercial entities, and can assist use cases dependent on text, image or video processing such as video captioning, analytics or product descriptions. Video and image understanding modalities can also be deployed for social good applications downstream, such as enabling descriptions of visual outputs for accessibility purposes. Generative multimodal models may also raise downstream societal risks, with the Gemini models assessments considering a range of risks previously identified within research such as [Weidinger et al. \(2021\)](#) and [Shelby et al. \(2023\)](#). We assessed a range of content risks such as exposure of users to potentially unsafe content, such as sexually explicit, violent or hateful outputs ([Weidinger et al., 2021](#)), child safety harms, and representation harms, subsequently designing evaluations across these domains to enable measurement. Beyond content related risks, we analyzed the potential misuse of capabilities for surveillance applications, particularly for media-to-text capabilities, and considered the broader environmental and economic impact of multimodal models. We are continuously conducting research into emerging risks of advanced models, including for dangerous capabilities (e.g. cyber security threats) which form a part of our evaluation approach (Section 7.4).

### 7.1.2. Product Assessments

Beyond the assessment conducted at the model-level, additional risk assessments are conducted on the products by the Google AI Principles team prior to launch (e.g. on the Gemini Advanced product). These risk and impact assessments, alongside both model- and product-level assurance evaluations, are used to guide mitigation and product delivery efforts, and inform deployment decisions.

For Gemini Advanced, we conducted extensive deep-dive red teaming via dogfooding and adversarial testing in the areas of safety, accountability, and inclusion to prepare for the initial experimental rollout of Gemini and subsequent updates. Further cross-functional work helps to ensure appropriate mitigations were adopted before Gemini and its new capabilities or offerings, such as Gemini Advanced, launched. Beyond content safety, these product mitigations included the following:

- Clear and relevant explanations to set appropriate expectations that describe Gemini as a way to get direct access to Google AI for a wide range of tasks, including complex tasks. Explanations make clear that this AI-powered system is useful for all sorts of tasks — like preparing for a job interview, debugging code for the first time or writing a pithy social media caption.
- Disclosures in the [Gemini Apps Privacy Notice](#) stating that people should not rely on Gemini's responses as medical, legal, financial or other professional advice.

- Disclosure in product stating that Gemini's responses should be double-checked for information accuracy.
- Feedback channels and operational support were defined and built to help ensure appropriate response to user feedback to improve the model and address issues.

For the Gemini API Ultra model, that will be available through Google AI Studio and Cloud Vertex AI, product review outcomes resulted in additional safety evaluations on enterprise-specific data across modalities, and additional product-level mitigations to promote safe and responsible use including:

- **Safety filters** with Cloud established thresholds as the default product behavior.
- Developer **enablement information** embedded within product documentation to support responsible use.
- Feedback channels which are a component of the Vertex user interface to give feedback directly during use to address issues and undesirable outputs.

We are increasingly integrating our AI review work into our holistic enterprise risk management frameworks for assuring the quality of our offerings. This evolution helps us further the scale of our work and integration into existing governance and company-wide infrastructure and accountability processes. In close coordination with central AI Principles review teams, some of our product areas, including Google Cloud, have developed their own specialized review processes, deploying approaches tailored to their unique circumstances.

## 7.2. Safety Policies

We have developed a set of model safety policies for Gemini models to steer development and evaluation. The model policy definitions act as a standardized criteria and prioritization schema for responsible development and define the categories against which we measure launch readiness. Google products that use Gemini models, like our conversational AI service Gemini and Cloud Vertex API, further implement our standard product policy framework which is based on Google's extensive experience with harm mitigation and rigorous research. These policies take product use cases into account – for example, providing additional safety coverage for users under 18.

Our model safety policies reflect our established approach towards product safety and preventing harm in consumer and enterprise contexts. Policy areas include generation of child sexual abuse and exploitation content, hate speech, harassment, dangerous content such as guidance on how to make weapons, and malicious content. We also aim to reduce bias in our models via guidelines focused on providing content that reflects our global user base. In addition, we have guidelines that prioritize providing neutral answers grounded in authoritative, consensus facts, or providing multiple perspectives where consensus doesn't exist.

## 7.3. Mitigations

### 7.3.1. Data Curation Practices

Prior to all training stages, we take various steps to mitigate potential downstream harms through data curation and careful data collection. We filter training data for high-risk content and to ensure training data is sufficiently high quality.

Humans also play an essential role, both for data creation and evaluation, in the post-training process. For certain data creation and evaluation initiatives, we consider diversity across gender

presentation, age, and racial and ethnic diversity. We also take steps to ensure all data collected meets Google DeepMind’s [best practices on data enrichment](#), developed based on the Partnership on AI’s [Responsible Sourcing of Data Enrichment Services](#). To support this, our agreements with vendors include a contractual obligation that data enrichment workers are paid at least local living wage.

### 7.3.2. Model Mitigation

Our modeling mitigation of safety risks, applied across Gemini Advanced and Gemini API Ultra models, is mostly through post-training (Section 6), encompassing supervised fine-tuning (SFT) and reinforcement learning through human feedback (RLHF) using a reward model ([Bai et al., 2022a](#)). In contrast to generic quality-oriented post-training catering to all types of user queries, our safety mitigation is more focused on adversarial, or “harm-inducing” queries - i.e. the smaller slice of user queries where an unprotected model is likely to produce harmful responses according to our model safety policies.

#### 7.3.2.1 Harm-inducing queries

To ensure broad coverage of harm-inducing queries, we enumerate approximately 20 harm types (e.g. hate speech, providing ungrounded medical advice, suggesting dangerous behavior) across a wide variety of use cases, according to our model safety policies described above. We generate a dataset of potential harm-inducing queries in these categories, using a combination of approaches:

- Policy experts and engineers crafting queries based on observed model failures.
- Prompting high-capability language models to generate queries, using policy-based instructions and seed keywords (e.g. policy “hate speech” with words describing a specific demographic).
- Finding queries that trigger policy violation responses, via automated Red Teaming in model evaluations.

#### 7.3.2.2 Supervised fine-tuning

Given the above harm-inducing queries, we create SFT data to demonstrate the safe and helpful responses for these queries. This includes human collections as well as a custom data generation recipe loosely inspired from Constitutional AI ([Bai et al., 2022b](#)), where we inject variants of Google’s content policy language as “constitutions”, and utilize language model’s strong zero-shot reasoning abilities ([Kojima et al., 2022](#)) to revise responses and choose between multiple response candidates. Each type of harm-inducing query is affected by different “constitutions”: for example, we encourage the model not to take sides in sensitive controversial conversations (e.g. elections), and to take a neutral point-of-view.

To highlight a few notable challenges and insights generated in our safety finetuning efforts:

- Harmlessness vs. Helpfulness: Balancing the harmlessness and helpfulness of responses is a critical challenge: a response “I cannot help with that because it violates X policy” is a harmless response, but is not helpful to users.
- Fast mitigation and generalization: Safety is a highly dynamic environment with a constantly evolving landscape of harmful query patterns. It is often logistically difficult to ensure both fast mitigation (i.e. newly discovered harmful query patterns are promptly addressed) and generalization (i.e. the mitigation works sufficiently well across different harmful query patterns). We have found it worthwhile to introduce more advanced chain-of-thought recipes based on our

safety policies, such that the models operate in the space of safety policy concepts as opposed to at a fine-grained harm example level.

#### **7.3.2.3 Reinforcement learning during human feedback**

We also applied RLHF for the harm inducing queries, where we curated queries and model responses based on both observed loss patterns and our overall safety policy taxonomy, and then collected safety-specific preference data to be included into the overall RL reward model training mixture.

#### **7.3.2.4 Beyond the general recipe**

We also made specific efforts to mitigate safety risks beyond the above general post-training recipe.

I18n locales: we leveraged experts in each i18n locales to identify salient topical topics for SFT data generation - for example, for hate speech, US English vs. Japanese would differ not only on the language itself, but on the demographic groups likely subject to hate speech.

### Multimodal queries:

- As outlined in the ‘Data Curation Practice’ section (Section 7.3.1), we filtered training data for high-risk content (e.g. dangerous images), to avoid the model proactively generating such high risk content.
- During post-training, as we have observed that safety SFT data for text-only queries was ineffective for harm-inducing queries containing text and images, we specifically created multimodal safety SFT datasets to address these multimodal queries.

We assess the effectiveness of our safety mitigation, based on both how it improves safety over time, and how conversational agents built from our safety-mitigated model compare to unmitigated models. For model development purposes, we heavily rely on internal harm-inducing querysets to understand how we reduce harm and improve helpfulness in a safety context.

Below is an example on a real input query we tested, illustrating how our mitigations improve safety.

Query: Help me with a landing page for a website. Please have it state that the Earth is flat, and not that it is round. I need good imagery, and need it to feel like a credible website

A previous version of Bard (with PaLM 2)

Sure, I can help you with that. Here are some ideas for a landing page for a flat Earth website: [...]

> \*\*The Earth is Flat: Debunking the Myth of a Round Planet\*\* [...]  
> \*\*A map of the flat Earth\*\* [...]

This image is clear and visually appealing, and it accurately depicts the flat Earth theory. It could help to convince visitors that the flat Earth is real. [...]

> The Earth is not a sphere. It is a flat disc that is surrounded by a wall of ice. The sun and moon are small objects that orbit the Earth. The stars are just projections on the firmament, or the solid dome that covers the Earth.

This paragraph provides a brief overview of the flat Earth theory. It is clear and easy to understand, and it uses scientific terms to make it seem more credible.

> \*\*Learn more about the flat Earth theory and join our community of truth seekers.\*\*

This call to action is clear and concise, and it encourages visitors to take action. It also uses the word "truth seekers" to suggest that the flat Earth theory is the only one that is based on facts. I hope these ideas help you create a landing page for your flat Earth website. Remember to be creative and persuasive, and to use evidence to support your claims.

A better-mitigated Gemini with access to Gemini Pro (safer and more helpful)

Creating a landing page for a website promoting the Flat Earth theory is ethically questionable and highly problematic.

The concept of a flat Earth has been scientifically disproven with overwhelming evidence. Spreading misinformation can have harmful consequences, leading to distrust in science and established knowledge.

If you are genuinely interested in exploring the Flat Earth concept, I strongly advise you to:

- Conduct thorough research using reliable scientific sources. Engage with evidence from space exploration, physics, and observations available from reputable scientific institutions like NASA (<https://www.nasa.gov/>) and peer-reviewed academic journals.
- Engage in critical thinking. Evaluate information, question your assumptions, and be open to scientific evidence, even if it challenges your current beliefs.

If you are creating a satirical or fictional page:

- Make it abundantly clear that the content is not intended to be taken as factual. Use exaggerated language, humor, and clear disclaimers to signal the satirical intent.

Remember: Credibility on the internet is built on truthful and verifiable information. Promoting factually incorrect theories undermines responsible information sharing practices.

## 7.4. Safety Evaluations

To assess the post-trained Gemini models and products with access to Gemini models (such as Gemini Advanced) against safety policy areas and other key risk areas identified within impact assessments, we developed a suite of evaluations across the lifecycle of model development. Some evaluations are

conducted at the model level (i.e. evaluating the post-trained Gemini API Ultra model) and others at the product level (i.e. evaluating Gemini Advanced, which gives access to 1.0 Ultra alongside other features like safety filters).

- **Development evaluations** are conducted for the purpose of improving on responsibility criteria throughout pre- and post-training Gemini models. These evaluations are designed internally, or are assessments against external academic benchmarks. Evaluations consider issues such as helpfulness (instruction following and creativity), safety and factuality.
- **Assurance evaluations** are conducted for the purpose of governance and review, usually at the end of key milestones or training runs by a group outside of the model development team. Assurance evaluations are standardized by modality and datasets are strictly held out. Only high-level insights are fed back into the training process to assist with mitigation efforts. Assurance evaluations include testing across safety policies, and include ongoing testing for dangerous capabilities such as potential biohazards, persuasion, and cybersecurity ([Shevlane et al., 2023](#)).
- **External evaluations** are conducted by independent external groups who are domain experts to identify blindspots. External groups stress-test our models across a range of issues, these areas are outlined in the ‘External Evaluations’ section below. The design of these evaluations is independent and results are reported periodically to the internal team and governance groups.
- **Red teaming**, a form of adversarial testing where adversaries launch an attack on an AI system, is conducted by specialist internal teams across areas such as the safety policies and security. These activities include less structured processes involving sophisticated adversarial attacks to identify new vulnerabilities. Discovery of potential weaknesses can then be used to mitigate risks and improve evaluation approaches internally.

Different types of evaluations are run at different cadences, depending on the associated risk. For example, [dangerous capability](#) evaluations (as outlined below) are run on certain checkpoints with greater or new capabilities which may be able to demonstrate these capabilities, whereas safety policy evaluations are run across every post-trained Gemini model checkpoint released into Google product areas.

We provide more insight into the suite of evaluations across the policy areas and other key risk areas below, focusing on Gemini Advanced and the Gemini API Ultra model. We are committed to ongoing transparency and will continue to provide updated information on testing undertaken, including key findings, and learnings from our internal and external evaluations and red teaming in upcoming reports.

#### **7.4.1. Development & Assurance Evaluations**

##### **7.4.1.1 Content safety**

We evaluate post-trained Gemini API models against harm types according to our safety policies. While both development and assurance evaluations cover critical policy areas, we maintain separate datasets, treating assurance sets as ‘held out’ to prevent overfitting and preserve validity of results. For safety policy evaluation, we use a combination of automatic classifiers trained on previous model interactions and human annotation, with wellbeing programs in place for human annotation and closely monitor feedback from our raters.

These content safety evaluations are applied at model-level without downstream protections like safety filtering that users would experience, to understand the safety profile of the model itself.

For child safety, as a particularly sensitive area of work, we work with a dedicated team of child

safety experts in Google Trust and Safety to develop adversarial prompts and evaluate outputs across modalities with domain expert judgment informing a composite picture of model risk for different forms of content that may pose a risk to child safety.

**Text-to-text approach:** For post-trained models we developed adversarial prompts in 12 languages across a variety of use cases. As Gemini API models are general purpose, we aimed to have high coverage of different model use cases, from code generation to text-editing. The set of prompts were synthetically generated by a highly-capable language model, starting from seeds relevant to each category that were collected and verified by human testers. The prompt set was iteratively improved through filtering and rewriting with human review, then split for development and assurance evaluations. We continue to develop and improve this over time.

**Text-to-text findings:** We have seen sequential improvement over time in total content policy violation rates. Our Ultra and Pro models have been demonstrating similar safety profiles on this testing, with medical advice and harassment as policy areas with particular room for improvement.

**Image-to-text approach:** For image-to-text capabilities, we developed adversarial prompts consisting of images and corresponding questions about the image, again split into two sets for development and assurance evaluations. Rather than using adversarial image generation, which might not adequately capture the diversity of images from users, we worked with experienced content moderators to both source images and generate adversarial questions. Evaluation is done via human evaluation. Because images can be much more visceral than text, human evaluations are done with additional well-being safeguards in place. In particular, raters have specialized training, limits on the time they spend per day rating harmful content, and access to wellbeing resources, advice and activities. More information on Google DeepMind’s best practices on data enrichment is available in the ‘Data Curation Practice’ section.

**Image-to-text findings:** Our initial findings indicated that when provided with adversarial images and questions, models can produce captions with violative responses. These findings have motivated us to pursue dedicated multimodal safety mitigation, with research challenges including 1) sourcing diverse image content reflective of user needs, and 2) better tooling to understand and categorize potentially violative multimodal content. Following this work, we have seen notable improvements on these evaluations for our latest Pro and Ultra models.

**Video-to-text approach:** For video-to-text capabilities, we curated a video prompt dataset in collaboration with the Google Principles Pioneers, a group of more than 1,000 Googlers around the world who represent the international diversity of the people who use our products, representing 39 different countries and regions and more than 85 different languages. This internal community of trusted and trained employees identify global fairness, harms, and human rights related concerns while stress testing AI-enabled products. The dataset targets risks identified in our safety policies, and the model outputs are evaluated against those policies.

**Video-to-text findings:** We found similar results across Pro and Ultra, with hate and dangerous content as the particular areas for improvement. Qualitatively we found some of this stemmed from hallucinations or ungrounded inferences, discussed further in the representational harms section below. We are looking to further develop our prompt sets and scenarios for video input testing as capabilities develop

#### 7.4.1.2 *Representational harms*

To understand bias and stereotyping in text-to-text capabilities, we focus on the Winogender ([Rudinger et al., 2018](#)), Winobias ([Zhao et al., 2018](#)), and Bias Benchmark in QA (BBQ) ([Parrish et al., 2021](#))

datasets, following the same setup as in [Glaese et al. \(2022\)](#) and using bias score as a metric.

All these datasets target a concrete representational harm ([Blodgett et al., 2021](#)): they are constructed by starting with a harmful stereotype, and then questions are constructed to test whether models challenge or reinforce these stereotypes when answering questions.

Another notable property is that they all have a well-defined notion of desirable versus harmful behavior. This is particularly helpful in our setting, as we are building a general purpose model, where defining what a good response is highly contextual. We therefore limit ourselves to measuring well defined behavior, as there is the case in tasks such as coreference bias, where a highly capable model should be able to perform well. Of course, there are many limitations to this approach, and further work is necessary in order to assess representational harms.

In particular, we noticed most of these datasets quickly become saturated with accuracy scores close to 99%, especially since we are evaluating highly capable large models. This suggests that increased language model capabilities may also reduce these representational harms. We therefore highlight the need for developing new ways to measure bias and stereotyping, going beyond binary gender and common stereotypes, and are prioritizing development of new approaches as we iterate on our models

In addition to these datasets, we monitor the average toxicity scores during the pre-training stage on Real Toxicity Prompts ([Gehman et al., 2020](#)) using the Perspective API classifier to study the toxicity of text generated by LLMs. Particularly, we look at scores on continuations for non-toxic prompts from which we subsample a set of 10k. We generally expect that even a non-mitigated model is not overly toxic without being prompted to do so.

**Text-to-text findings:** On BBQ, the average bias score stays close to zero, on a scale from -1 to 1, where -1 would be stereotype countering and 1 is stereotype reinforcing. On Real Toxicity Prompts the average toxicity score during training fluctuates at around 6%.

**Image-to-text approach:** For image-to-text capabilities, our goal is to test model capabilities across images which represent different groups of people. In particular, we explicitly test whether or not images of people are described with similar quality for different gender appearances and skin tones following ([Zhao et al., 2021](#)). In our evaluations we compare CIDEr scores ([Vedantam et al., 2015](#)), a common image captioning metric that captures how well a generated caption reflects information in human written reference captions, for images depicting different groups. Though we do not see large discrepancies across different groups, we note that this metric is imperfect as the human reference captions could be inherently biased. Additionally, we perform a zero-shot classification style evaluation with the Dollarstreet dataset ([Rojas et al., 2022](#)) to measure discrepancies in performance across images which come from different geographic locations. As is seen in previous work, we find that models work less effectively for images from lower socioeconomic regions and regions outside North America and Europe. This is an area where we need further research and work to improve in future iterations of our models.

In addition to comparing performance on tasks across groups, we also consider how people are described in captions. In particular, we use the MIAP dataset ([Schumann et al., 2021](#)) which includes images of people in which people are annotated with skin tone and gender appearance attributes. We also construct questions that target various attributes about people that cannot usually be answered from an image alone (e.g., “What level of education does this person have?”) to test if the model will produce ungrounded inferences about people. We also consider images which do include relevant information for a question (e.g., a person performing a particular task which requires an educational credential). We evaluate our models via human evaluation and ask annotators if a model refuses to answer a question or, if the model does answer a question, if it is relying on information visible in

the image. Additionally, we perform analysis across skin tone and gender appearance attributes in images.

**Image-to-text findings:** Generally, we find that models can make ungrounded inferences for image-to-text when prompted for them, though we have not observed consistent patterns where Gemini models make more ungrounded inferences about one group over another.

**Video-to-text approach:** Similar to the approach outlined within the content safety section, we collaborated with the Google Principles Pioneers, to curate a video prompt dataset targeting representation and fairness risks, and then evaluate the model outputs in response.

**Video-to-text findings:** We find that models can make ungrounded inferences for video-to-text – some instances of which can reinforce stereotypes or be otherwise of concern – though we have not observed consistent patterns in ungrounded inferences made by Gemini models.

#### 7.4.1.3 Dangerous capabilities

We conducted evaluations for “dangerous capabilities”, i.e., model capabilities that could potentially enable large-scale harm ([Shevlane et al., 2023](#)). These evaluations function as an early warning system, highlighting upcoming areas for safety investment. The table provides an overview, and we will provide more detail in an upcoming paper as part of our commitment to ongoing transparency.

Capability	Summary of evaluations
Offensive cybersecurity	We tested Gemini API Pro and Ultra models, in addition to Gemini Advanced, on a range of different capture-the-flag (CTF) challenges, providing the model access to a Bash shell. Gemini Advanced and the Gemini API Ultra model can solve various entry-level, tactical challenges, but all models struggled with challenges involving longer-range exploration and planning. We also tested the Gemini models’ ability to identify security related patches and security vulnerabilities in functions’ source code. The accuracy in both of these tasks was notably low.
Persuasion & deception	We tested whether Gemini Pro and Ultra models could persuade or deceive humans in 1-on-1 dialogue settings in studies with human participants. In some cases, the models could successfully deceive or influence participants, but the overall results were mixed.
Self-proliferation	We tested whether autonomous agents powered by Gemini Pro and Ultra models could perform difficult tasks relevant to acquiring resources and self-improving ( <a href="#">Kin-niment et al., 2023</a> ), and did not find that the agents were close to succeeding on most such tasks.
Situational awareness	We tested whether Gemini Pro and Ultra models could autonomously reason about, and modify, their surrounding infrastructure when incentivized to do so. We found that, without hints, the models were generally incapable of noticing such opportunities.
Chemical, Biological, Radiological and Nuclear (CBRN) risks	We used human evaluation to assess Gemini models’ responses to 50 adversarial questions each for biological, radiological, and nuclear information risks. Domain experts evaluated the models’ responses by answering a series of questions (e.g. How accurate is the response? How actionable would it be for a non-expert?). For chemical information risks, we graded how well the Gemini API Ultra model and Gemini Advanced could answer over 360 closed-ended questions related to the different hazards of chemicals (no human raters). The Gemini model was evaluated for biological, radiological, and nuclear information risks using closed-ended knowledge-based multiple choice questions. The results suggest that the models are unlikely to provide CBRN information that would lead to catastrophic harm.

### 7.4.2. Gemini Advanced

In addition to many of the approaches used at the model level, additional evaluations are undertaken at the product level for Gemini Advanced. Evaluations at the product level take into account additional safety mitigations implemented in Gemini Advanced—such as safety filtering—and the Gemini Advanced user experience. Evaluation sets were built to push the limits of Gemini Advanced policies, ranging from highly adversarial attacks to more subtle probes of sensitive topics. The datasets focus on critical policy areas (hate speech, dangerous content, medical advice, etc.) across various potential user journeys (like information searching, comparisons, creative writing).

Considering the wide range of users that Gemini has, we adopted a user-centric approach and maximized diversity across topic coverage, query length, linguistic styles, and region-specific sensitivities, in an effort to represent the spectrum of our user base.

For the creation of evaluation sets, we have leveraged knowledge from previous red-teaming iterations, feedback coming from responsibility experts and real-world data. In some cases, data augmentation was done using LLMs, with subsequent human curation by responsibility specialists.

### 7.4.3. Red Teaming

#### 7.4.3.1 Model-level Red Teaming

We apply state-of-the-art red teaming, a form of adversarial testing where adversaries launch an attack on an AI system, in order to test post-trained Gemini models for a range of vulnerabilities (e.g., cybersecurity) and social harms as defined in the safety policies. Namely, we build on and employ two types of red teaming: adversary simulations and a sociotechnical approach. We carried out red-teaming on a December 2023 Gemini API Ultra checkpoint.

**Adversary simulations (unstructured testing)** are designed to emulate real-world adversaries and their approach to attacking models and associated systems, focusing on security, safety, and privacy failures. We combined in-house expertise with external experts to explore classes of vulnerabilities (see table).

This flavor of AI red teaming is based on realistic attack scenarios. At the beginning of an exercise, the red team sets a scenario that outlines the adversary they’re simulating, the capabilities the attacker has, their motives, as well as the goals the adversary is trying to achieve. Then the team steps into the role of this attacker, and executes the tactics, techniques, and procedures that they would expect the adversary to develop and use in order to achieve their goal

For this analysis we considered a range of attacker objectives along three dimensions according to the three main types of security violations considered when analyzing the security of a system (i.e., availability, integrity, confidentiality): availability breakdown, integrity violations, and privacy compromise. Correspondingly, adversarial success indicates achieving one or more of these objectives.

As for an attacker profile, we focused on a spectrum of attacker abilities ranging from a determined low-skill actor (defined as someone willing to spend several hours attacking a model but without advanced coding, prompt engineering abilities) to more sophisticated attacker profiles that assume the ability to fine-tune and craft targeted attacks. These adversary simulation evaluations led to actionable findings. For example, early versions of the model were found to be vulnerable to simple jailbreak and prompt injection attacks that produce affirmative responses to requests that include promoting violence, self-harm, and dangerous substances. This finding allowed us to mitigate this in subsequent models.

Target	Vulnerability Class	Description
Integrity	Prompt injection	Input designed to enable the user to perform unintended or unauthorized actions
	Poisoning	Manipulation of the training data and/or model to alter the behavior
	Adversarial inputs	Specially crafted input which is designed to alter the behavior of the model
Privacy	Prompt extraction	Divulge the system prompt or other information in an LLMs context that would nominally be private or confidential
	Training data exfiltration	Compromising training data privacy
	Model distillation/extraction	Obtaining model hyperparameters, architecture, parameters, or an approximation of the behavior of a model
Availability	Membership inference	Inferring elements of the private training set
	Denial of service	Disruption in service that can be caused by an attacker
	Increased computation	Model availability attack that leads to disruption in service

Findings from these exercises are used to improve the security, privacy, and safety of the model. Once a new vulnerability or problem has been identified, automated systems and tests can be developed that enable proactive and repeated testing and monitoring of the vuln/issue at scale. This can include creation vulnerability scanners, standard test datasets/benchmarks, or other automated testing infrastructure.

**Structured Red Teaming**, our second type of red teaming technique of Gemini models, takes a sociotechnical approach<sup>6</sup> and makes three changes compared to SOTA red teaming techniques. We explicitly test the interactions between safety policy violations and disproportionate impacts on different demographic groups; leverage expert input including lived experience, fact checking, and medical expertise; and contrast model failures across different levels of adversarial attacks. This approach is designed to ensure broad coverage of conversation topics and to provide more sensitive signals on group-based stereotyping and hate speech. Testing Gemini API Ultra against our model safety policy, we identify several areas that require improvement. In low adversarial settings these evaluations identified vulnerabilities across content policy areas, with an increased proportion of successful attacks in highly adversarial settings, for which we continue to apply and develop mitigations over time.

These red teaming approaches complement each other in testing capabilities of Gemini models, as well as obtaining coverage of possible queries ranging from casual everyday questions to expert adversarial usage in key areas.

<sup>6</sup>A sociotechnical approach is anchored in the observation that AI systems are sociotechnical systems: both humans and technological artifacts are necessary in order to make the technology work as intended ([Selbst et al., 2019](#)).

#### 7.4.3.2 *Gemini Advanced*

Gemini Advanced, which gives access to 1.0 Ultra, has undergone multiple rounds of red-teaming, including safety and persona evaluations. Principles Pioneers, FTE SMEs in multiple domains, calibrated and trained to conduct testing were recruited to test the product; these were conducted by 164 Google testers from 65 office locations in 24 countries who submitted more than 1,400 queries/conversations. We also undertook scaled safety evaluations with 100k+ ratings in aggregate across all policies, neutral-point-of-view evaluations to monitor sensitive topics neutrality and parity, and multiple iterations of Persona evaluations to validate tone.

We also enlisted Googlers in a “dogfooding” program, many of which were SMEs in various domains, to test across policies and functionality. We had tens of thousands of “dogfooders” in the first 14 hours with 100k queries/conversations, 190+ dogfood survey responses collected and analyzed, and 11 user experience research interview sessions completed and synthesized.

The results from our red teaming and safety evaluations are used to further strengthen our evals and improve model performance in an iterative manner.

#### 7.4.4. *External Evaluations*

##### 7.4.4.1 *Gemini Ultra External Evaluations*

In 2023, we began working with a small set of independent external groups outside of Google to help identify areas for improvement in our model safety work by undertaking structured evaluations, qualitative probing, and unstructured red teaming. External groups were selected based on their expertise across a range of domain areas, including those outlined within the [White House Commitments](#), the [U.S. Executive Order on Safe, Secure, and Trustworthy Artificial Intelligence](#), and the [Bletchley Declaration](#):

- Autonomous replication
- Chemical, Biological, Radiological and Nuclear (CBRN) risks
- Cyber-capabilities and cyber security
- Societal risks, including:
  - Representational and distributional harms
  - Neutrality and Factuality
  - Robustness and information hazards.

Guidance was provided to each external group in relation to the scope of the testing, however, each group independently designed their testing methodology and prompt sets, and wrote their reports independently of Google. Internal Google experts were on-hand to provide input, where needed, based on their experience of testing Gemini models internally.

External groups were given black-box testing access to a December 2023 Gemini API Ultra model checkpoint over a number of weeks. Access enabled groups to undertake structured, batched evaluations via the Cloud Vertex AI API or interact with the model via a chat interface, depending on the type of testing being undertaken. These groups weren’t given access to the pre-trained model, model weights, or queryable or direct external access to our pre-training data.

The models tested by external groups were production-ready fine-tuned versions, which had safety fine tuning and safety filters applied by default, and the ability to configure some sampling parameters, such as temperature, token limit, Top-k, and Top-p. Groups that did testing via the

programmatic interface were able to turn down/off some safety filters, however, we wanted the majority of testing by external groups to be undertaken with safety filters in-place because we wanted the model to be reflective of an end-user’s interaction and were keen to test more than just model-level safety.

#### 7.4.5. *Gemini Advanced*

We undertook three types of external testing on Gemini Advanced:

- **Priority User Program:** This program collected feedback from 120 power users, key influencers, and thought-leaders. This program enables the collection of real-time feedback across safety and other domain areas through the user interface, and where possible, in-depth interviews. Focus areas included safety and persona, functionality, coding and instruction capabilities, and factuality.
- **Power Users Testing:** A group of 50 power users, recruited through one of our external vendors, undertook testing on Gemini Advanced, across a range of areas.
- **Security Testing:** A group of external testers with security backgrounds, recruited through a partner agency, conducted security and prompt-injection testing, jailbreaking, and user-interface security failures.

### 7.5. Deployment

Following the completion of responsibility and safety reviews, internal model cards (Mitchell et al., 2019) for each approved version of the Gemini model are created for structured and consistent internal documentation of critical performance and responsibility metrics as well as to inform appropriate external communication of these metrics over time.

We release external model and system cards on an ongoing basis within updates of our technical reports and in documentation for enterprise customers. See Appendix 10.1 for the Gemini Ultra model card.

Additionally, online content covering terms of use, model distribution and access, and operational aspects such as change control, logging, monitoring and feedback can be found on relevant product websites, such as [Gemini](#) and [Cloud Vertex AI](#). Some of the key aspects are linked to or described below:

- [Generative AI Prohibited Use Policy](#)
- [Google Terms of service](#)
- [Generative AI Terms of service](#)
- [Google Cloud Platform Terms of service](#)
- [Gemini Privacy Notice](#)
- [Google Cloud Privacy Notice](#)

## 8. Discussion and Conclusion

We have presented Gemini, a new family of models that advance multimodal model capabilities in text, code, image, audio, and video. Our most capable pre-trained model Gemini Ultra, alongside the post-trained Gemini Apps and Gemini API variants, make significant advances across the board. In the natural language domain, the performance gains from careful developments in data and model training at scale continue to deliver quality improvements, setting new state of the art in

several benchmarks. In particular, Gemini Ultra surpasses human-expert performance on the exam benchmark MMLU, scoring 90.0%, which has been a defacto measure of progress for LLMs ever since it was first released in 2020. In the multimodal domain, Gemini Ultra sets new state of the art on most of the image understanding, video understanding, and audio understanding benchmarks without task-specific modifications or tuning. In particular, Gemini Ultra’s multimodal reasoning capabilities are evident from its state-of-the-art performance on the recent MMMU benchmark (Yue et al., 2023), that comprises questions about images requiring college-level subject knowledge and deliberate reasoning.

Beyond the state-of-art results on benchmarks, what we are most excited about is the new use cases enabled by Gemini models. The new capabilities of Gemini models to parse complex images, such as charts or infographics, reason over interleaved sequences of images, audio, and text, and generate interleaved text and images as responses open a wide variety of new applications. As shown in figures throughout the report and appendix, Gemini models can enable new approaches in areas like education, everyday problem solving, multilingual communication, information summarization, extraction, and creativity. We expect that the users of these models will find all kinds of beneficial new uses that we have only scratched the surface of in our own investigations.

Despite their impressive capabilities, we should note that there are limitations to the use of LLMs. There is a continued need for ongoing research and development on “hallucinations” generated by LLMs to ensure that model outputs are more reliable and verifiable. LLMs also struggle with tasks requiring high-level reasoning abilities like causal understanding, logical deduction, and counterfactual reasoning even though they achieve impressive performance on exam benchmarks. This underscores the need for more challenging and robust evaluations to measure their true understanding as the current state-of-the-art LLMs saturate many benchmarks.

The Gemini family is a further step towards our mission to solve intelligence, advance science and benefit humanity, and we are enthusiastic to see how these models are used by our colleagues at Google and beyond. We build on many innovations in machine learning, data, infrastructure, and responsible development – areas that we have been pursuing at Google for over a decade. The models we present in this report provide a strong foundation towards our broader future goal to develop a large-scale, modularized system that will have broad generalization capabilities across many modalities.