August 22, 2023

# GPT-3.5 Turbo fine-tuning and API updates

Developers can now bring their own data to customize GPT-3.5 Turbo for their use cases.
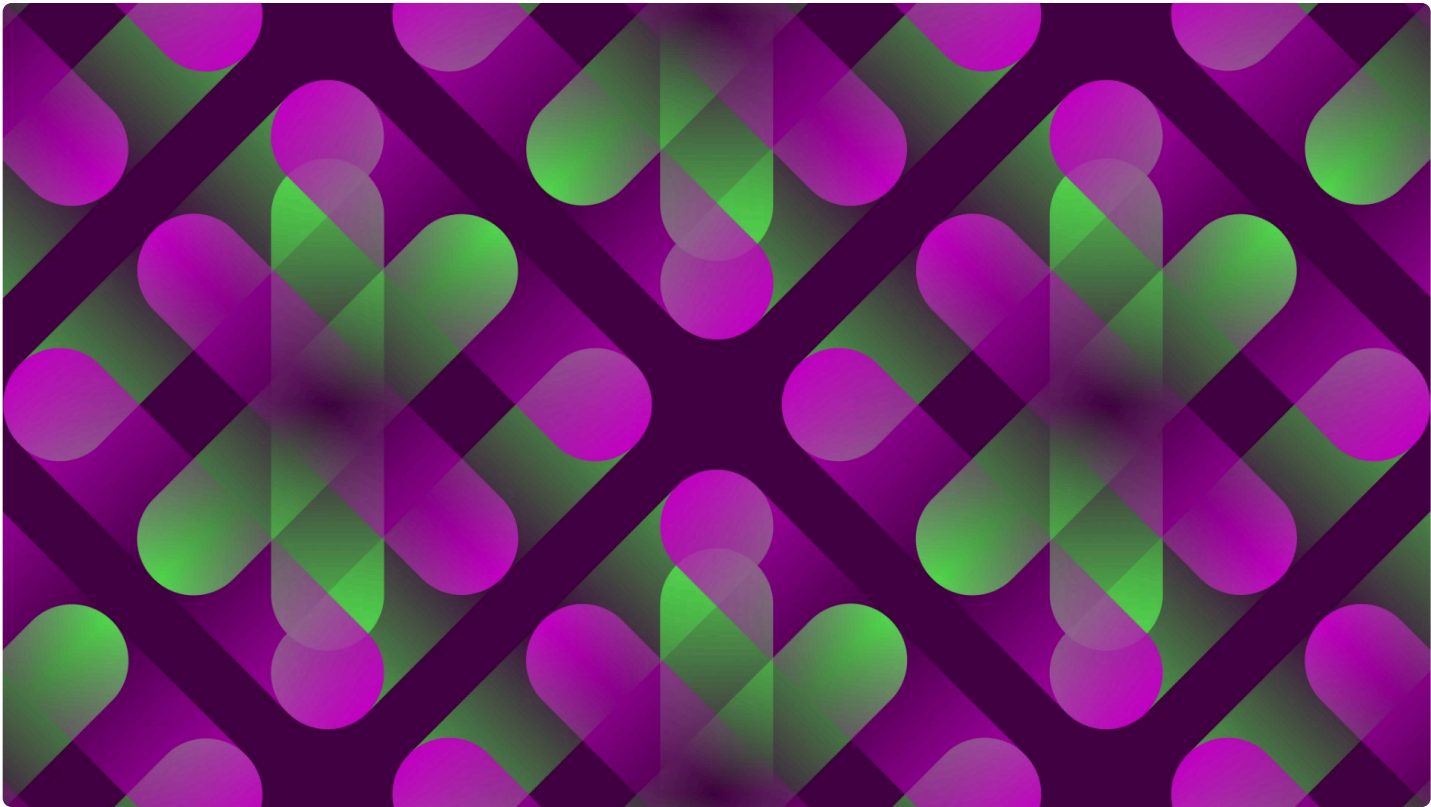


Illustration: Ruby Chen

Fine-tuning for GPT-3.5 Turbo is now available, with fine-tuning for GPT-4 coming this fall. This update gives developers the ability to customize models that perform better for their use cases and run these custom models at scale. Early tests have shown a fine-tuned version of GPT-3.5 Turbo can match, or even outperform, base GPT-4-level capabilities on certain narrow tasks. As with all our APIs, data sent in and out of the fine-tuning API is owned by the customer and is not used by OpenAI, or any other organization, to train other models.

## Fine-tuning use cases

Since the release of GPT-3.5 Turbo, developers and businesses have asked for the ability to customize the model to create unique and differentiated experiences for their users. With this launch, developers can now run supervised fine-tuning to make this model perform better for their use cases.

In our private beta, fine-tuning customers have been able to meaningfully improve model performance across common use cases, such as:

- **Improved steerability**: Fine-tuning allows businesses to make the model follow instructions better, such as making outputs terse or

German when prompted to use that language.

- **Reliable output formatting:** Fine-tuning improves the model's ability to consistently format responses—a crucial aspect for applications demanding a specific response format, such as code completion or composing API calls. A developer can use fine-tuning to more reliably convert user prompts into high-quality JSON snippets that can be used with their own systems.

- **Custom tone:** Fine-tuning is a great way to hone the qualitative feel of the model output, such as its tone, so it better fits the voice of businesses' brands. A business with a recognizable brand voice can use fine-tuning for the model to be more consistent with their tone.

In addition to increased performance, fine-tuning also enables businesses to **shorten their prompts** while ensuring similar performance.  Fine-tuning with GPT-3.5-Turbo can also handle 4k tokens—double our previous fine-tuned models. Early testers have

reduced prompt size by up to 90% by fine-tuning instructions into the model itself, speeding up each API call and cutting costs.

Fine-tuning is most powerful when combined with <u>other techniques</u> such as prompt engineering, information retrieval, and function calling. Check out our <u>fine-tuning guide</u> to learn more. Support for fine-tuning with function calling and `gpt-3.5-turbo-16k` will be coming later this fall.

## Fine-tuning steps

| Step 1 | |
|---|---|
| **Prepare your data** | ‹› |

| Step 2 | |
|---|---|
| **Upload files** | ‹› |

| Step 3 | |
|---|---|
| **Create a fine-tuning job** | ‹› |

Once a model finishes the fine-tuning process, it is available to be used in production right away and has the same shared rate limits as the underlying model.

| Step 4 | |
|---|---|
| **Use a fine-tuned model** | ‹› |

We will also be debuting a fine-tuning UI in the near future, which will give developers easier access to information about ongoing fine-tuning jobs, completed model snapshots, and more.

## Safety

It is very important to us that the deployment of fine-tuning is safe. To preserve the default model's safety features through the fine-tuning process, fine-tuning training data is passed through our Moderation API and a GPT-4 powered moderation system to detect unsafe training data that conflict with our safety standards.

## Pricing

Fine-tuning costs are broken down into two buckets: the initial training cost and usage cost:

- Training: $0.008 / 1K Tokens

- Usage input: $0.012 / 1K Tokens

- Usage output: $0.016 / 1K Tokens

For example, a `gpt-3.5-turbo` fine-tuning job with a training file of 100,000 tokens that is trained for 3 epochs would have an expected cost of $2.40.

## Updated GPT-3 models

In July, <u>we announced</u> that the original GPT-3 base models ( `ada` , `babbage` , `curie` , and `davinci` ) would be turned off on January 4th, 2024. Today, we are making `babbage-002` and `davinci-002` available as replacements for these models, either as base or fine-tuned models. Customers can access those models by querying the <u>Completions API</u>.

These models can be fine-tuned with our new API endpoint `/v1/fine_tuning/jobs` . This new endpoint offers pagination and more extensibility to support the future evolution of the fine-tuning API. Transitioning from `/v1/fine-tunes` to the updated endpoint is straightforward and more details can be found in our new <u>fine-tuning guide</u>. This deprecates the old `/v1/fine-tunes` endpoint, which will be turned off on January 4th, 2024.

Pricing for base and fine-tuned GPT-3 models is as follows:

| Model | Base models | | Fine-tuned models | | |
|---|---|---|---|---|---|
| | Input tokens | Output tokens | Training | Input tokens | Output tokens |
| babbage-002 | $0.0004 / 1K tokens | $0.0004 / 1K tokens | $0.0004 / 1K tokens | $0.0016 / 1K tokens | $0.0016 / 1K tokens |
| davinci-002 | $0.002 / 1K tokens | $0.002 / 1K tokens | $0.006 / 1K tokens | $0.012 / 1K tokens | $0.012 / 1K tokens |

## Authors

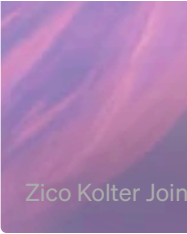Andrew Peng, Michael Wu, John Allard, Logan Kilpatrick, Steven Heidel

## Acknowledgments

Andrea Vallone, Arvind Neelakantan, Cindy Yong, Colin Jarvis, Denny Jin, Florencia Leoni Aleman, Henry Head, Ilan Bigio, Jeff Harris, Jessica Shieh, Juston Forte, Kim Malfacini, Lauren Workman, Lilian Weng, Olivier Godement, Sherwin Wu, Shyamal Anadkat, Vik Goel, Yuchen He

## Related articles                                                                View all Company >

Aug 20, 2024

Aug 8, 2024

OpenAI partners with Condé Nast

Zico Kolter Join

Our research

Overview

Index

Latest advancements

GPT-4

GPT-4o mini

DALL·E 3

Sora

ChatGPT

For Everyone

For Teams

For Enterprises

ChatGPT login ↗
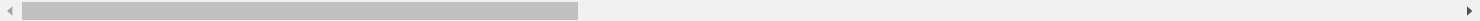
Download

API

Platform overview

Pricing

Documentation ↗

API login ↗

Explore more

OpenAI for business

Safety overview

Safety overview

Safety standards

Teams

Safety Systems

Preparedness

Superalignment

Company

About us

News

Our Charter

Security

Residency

Careers

Terms & policies

Terms of use

Privacy policy

Brand guidelines

Other policies

OpenAI © 2015–2024