

# **Fine-tuning**

Learn how to customize a model for your application.

### Introduction

Fine-tuning lets you get more out of the models available through the API by providing:

Higher quality results than prompting

Ability to train on more examples than can fit in a prompt

Token savings due to shorter prompts

Lower latency requests

OpenAl's text generation models have been pre-trained on a vast amount of text. To use the models effectively, we include instructions and sometimes several examples in a prompt. Using demonstrations to show how to perform a task is often called "few-shot learning."

Fine-tuning improves on few-shot learning by training on many more examples than can fit in the prompt, letting you achieve better results on a wide number of tasks. **Once a model has been fine-tuned, you won't need to provide as many examples in the prompt.** This saves costs and enables lower-latency requests.

At a high level, fine-tuning involves the following steps:

- Prepare and upload training data
- 2 Train a new fine-tuned model
- 3 Evaluate results and go back to step 1 if needed
- 4 Use your fine-tuned model

Visit our pricing page to learn more about how fine-tuned model training and usage are billed.

#### What models can be fine-tuned?

i Fine-tuning for GPT-4 is in an experimental access program - eligible users can request access in the fine-tuning UI when creating a new fine-tuning job.

Fine-tuning is currently available for the following models:



```
babbage-002
davinci-002
gpt-4-0613 (experimental)
```

You can also fine-tune a fine-tuned model which is useful if you acquire additional data and don't want to repeat the previous training steps.

We expect gpt-3.5-turbo to be the right model for most users in terms of results and ease of use.

# When to use fine-tuning

Fine-tuning OpenAI text generation models can make them better for specific applications, but it requires a careful investment of time and effort. We recommend first attempting to get good results with prompt engineering, prompt chaining (breaking complex tasks into multiple prompts), and function calling, with the key reasons being:

There are many tasks at which our models may not initially appear to perform well, but results can be improved with the right prompts - thus fine-tuning may not be necessary

Iterating over prompts and other tactics has a much faster feedback loop than iterating with fine-tuning, which requires creating datasets and running training jobs

In cases where fine-tuning is still necessary, initial prompt engineering work is not wasted - we typically see best results when using a good prompt in the fine-tuning data (or combining prompt chaining / tool use with fine-tuning)

Our prompt engineering guide provides a background on some of the most effective strategies and tactics for getting better performance without fine-tuning. You may find it helpful to iterate quickly on prompts in our playground.

#### Common use cases

Some common use cases where fine-tuning can improve results:

Setting the style, tone, format, or other qualitative aspects

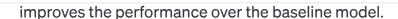
Improving reliability at producing a desired output

Correcting failures to follow complex prompts

Handling many edge cases in specific ways

Performing a new skill or task that's hard to articulate in a prompt





Another scenario where fine-tuning is effective is in reducing costs and / or latency, by replacing GPT-4 or by utilizing shorter prompts, without sacrificing quality. If you can achieve good results with GPT-4, you can often reach similar quality with a fine-tuned | gpt-3.5-turbo | model by finetuning on the GPT-4 completions, possibly with a shortened instruction prompt.

# **Preparing your dataset**

Once you have determined that fine-tuning is the right solution (i.e. you've optimized your prompt as far as it can take you and identified problems that the model still has), you'll need to prepare data for training the model. You should create a diverse set of demonstration conversations that are similar to the conversations you will ask the model to respond to at inference time in production.

Each example in the dataset should be a conversation in the same format as our Chat Completions API, specifically a list of messages where each message has a role, content, and optional name. At least some of the training examples should directly target cases where the prompted model is not behaving as desired, and the provided assistant messages in the data should be the ideal responses you want the model to provide.

# **Example format**

In this example, our goal is to create a chatbot that occasionally gives sarcastic responses, these are three training examples (conversations) we could create for a dataset:

```
1 {"messages": [{"role": "system", "content": "Marv is a factual cha
2 {"messages": [{"role": "system", "content": "Marv is a factual challet the
3 {"messages": [{"role": "system", "content": "Marv is a factual chatbot the
```

The conversational chat format is required to fine-tune | gpt-3.5-turbo |. For | babbage-002 | and davinci-002, you can follow the prompt completion pair format as shown below.

```
1 {"prompt": "prompt text>", "completion": "<ideal generated text>"
2 {"prompt": "prompt text>", "completion": "<ideal generated text>",
3 {"prompt": "rompt text>", "completion": "<ideal generated text>"}
```



We generally recommend taking the set of instructions and prompts that you found worked best for the model prior to fine-tuning, and including them in every training example. This should let you reach the best and most general results, especially if you have relatively few (e.g. under a hundred) training examples.

If you would like to shorten the instructions or prompts that are repeated in every example to save costs, keep in mind that the model will likely behave as if those instructions were included, and it may be hard to get the model to ignore those "baked-in" instructions at inference time.

It may take more training examples to arrive at good results, as the model has to learn entirely through demonstration and without guided instructions.

# **Example count recommendations**

To fine-tune a model, you are required to provide at least 10 examples. We typically see clear improvements from fine-tuning on 50 to 100 training examples with | gpt-3.5-turbo | but the right number varies greatly based on the exact use case.

We recommend starting with 50 well-crafted demonstrations and seeing if the model shows signs of improvement after fine-tuning. In some cases that may be sufficient, but even if the model is not yet production quality, clear improvements are a good sign that providing more data will continue to improve the model. No improvement suggests that you may need to rethink how to set up the task for the model or restructure the data before scaling beyond a limited example set.

# Train and test splits

After collecting the initial dataset, we recommend splitting it into a training and test portion. When submitting a fine-tuning job with both training and test files, we will provide statistics on both during the course of training. These statistics will be your initial signal of how much the model is improving. Additionally, constructing a test set early on will be useful in making sure you are able to evaluate the model after training, by generating samples on the test set.

# **Token limits**

Token limits depend on the model you select. For gpt-3.5-turbo-1106, the maximum context length is 16,385 so each training example is also limited to 16,385 tokens. For gpt-3.5-turbo-0613, each training example is limited to 4,096 tokens. Examples longer than the default will be truncated to the maximum context length which removes tokens from the end of the training



You can compute token counts using our counting tokens notebook from the OpenAl cookbook.

#### **Estimate costs**

Please refer to the pricing page for details on cost per 1k input and output tokens (we do to charge for tokens that are part of the validation data). To estimate the costs for a specific fine-tuning job, use the following formula:

base cost per 1k tokens \* number of tokens in the input file \* number of epochs trained

For a training file with 100,000 tokens trained over 3 epochs, the expected cost would be ~\$2.40 USD.

# **Check data formatting**

Once you have compiled a dataset and before you create a fine-tuning job, it is important to check the data formatting. To do this, we created a simple Python script which you can use to find potential errors, review token counts, and estimate the cost of a fine-tuning job.

#### Fine-tuning data format validation

Learn about fine-tuning data formatting

# Upload a training file

Once you have the data validated, the file needs to be uploaded using the Files API in order to be used with a fine-tuning jobs:

```
1 from openai import OpenAI
2 client = OpenAI()
3
4 client.files.create(
5 file=open("mydata.jsonl", "rb"),
6 purpose="fine-tune"
7 )
```



The maximum file upload size is 1 GB, though we do not suggest fine-tuning with that amount of data since you are unlikely to need that large of an amount to see improvements.

## Create a fine-tuned model

1 from openai import OpenAI

4 client.fine\_tuning.jobs.create(
5 training\_file="file-abc123",
6 model="gpt-3.5-turbo"

2 client = OpenAI()

3

6 7 )

After ensuring you have the right amount and structure for your dataset, and have uploaded the file, the next step is to create a fine-tuning job. We support creating fine-tuning jobs via the fine-tuning UI or programmatically.

To start a fine-tuning job using the OpenAI SDK:

```
python ∨ ☐ Copy
```

In this example, <code>model</code> is the name of the model you want to fine-tune (<code>gpt-3.5-turbo</code>), <code>babbage-002</code>, <code>davinci-002</code>, or an existing fine-tuned model) and <code>training\_file</code> is the file ID that was returned when the training file was uploaded to the OpenAl API. You can customize your fine-tuned model's name using the <code>suffix parameter</code>.

To set additional fine-tuning parameters like the validation\_file or hyperparameters, please refer to the API specification for fine-tuning.

After you've started a fine-tuning job, it may take some time to complete. Your job may be queued behind other jobs in our system, and training a model can take minutes or hours depending on the model and dataset size. After the model training is completed, the user who created the fine-tuning job will receive an email confirmation.

In addition to creating a fine-tuning job, you can also list existing jobs, retrieve the status of a job, or cancel a job.



```
2
    # List 10 fine-tuning jobs
3
    client.fine_tuning.jobs.list(limit=10)
4
5
    # Retrieve the state of a fine-tune
6
    client.fine_tuning.jobs.retrieve("ftjob-abc123")
7
8
    # Cancel a job
9
    client.fine_tuning.jobs.cancel("ftjob-abc123")
10
11
    # List up to 10 events from a fine-tuning job
12
    client.fine_tuning.jobs.list_events(fine_tuning_job_id="ftjob-abc123", limit=10)
13
14
    # Delete a fine-tuned model (must be an owner of the org the model was created in)
15
    client.models.delete("ft:gpt-3.5-turbo:acemeco:suffix:abc123")
16
17
```

# Use a fine-tuned model

When a job has succeeded, you will see the fine\_tuned\_model field populated with the name of the model when you retrieve the job details. You may now specify this model as a parameter to in the Chat Completions (for gpt-3.5-turbo) or legacy Completions API (for babbage-002 and davinci-002), and make requests to it using the Playground.

After your job is completed, the model should be available right away for inference use. In some cases, it may take several minutes for your model to become ready to handle requests. If requests to your model time out or the model name cannot be found, it is likely because your model is still being loaded. If this happens, try again in a few minutes.

```
1
    from openai import OpenAI
    client = OpenAI()
2
3
4
    response = client.chat.completions.create(
5
      model="ft:gpt-3.5-turbo:my-org:custom_suffix:id",
6
      messages=[
7
        {"role": "system", "content": "You are a helpful assistant."},
8
        {"role": "user", "content": "Hello!"}
      ٦
```



You can start making requests by passing the model name as shown above and in our GPT guide.

# **Analyzing your fine-tuned model**

We provide the following training metrics computed over the course of training: training loss, training token accuracy, test loss, and test token accuracy. These statistics are meant to provide a sanity check that training went smoothly (loss should decrease, token accuracy should increase). While an active fine-tuning jobs is running, you can view an event object which contains some useful metrics:

```
1
   {
        "object": "fine_tuning.job.event",
2
3
        "id": "ftevent-abc-123",
        "created_at": 1693582679,
4
5
        "level": "info",
6
        "message": "Step 100/100: training loss=0.00",
7
        "data": {
            "step": 100.
8
9
            "train_loss": 1.805623287509661e-5,
10
            "train_mean_token_accuracy": 1.0
11
        },
        "type": "metrics"
12
13 }
```

After a fine-tuning job has finished, you can also see metrics around how the training process went by querying a fine-tuning job, extracting a file ID from the result\_files, and then retrieving that files content. Each results CSV file has the following columns: step, train\_loss, train\_accuracy, valid\_loss, and valid\_mean\_token\_accuracy.

```
1 step,train_loss,train_accuracy,valid_loss,valid_mean_token_accurac
2 1,1.52347,0.0,,
3 2,0.57719,0.0,,
4 3,3.63525,0.0,,
5 4,1.72257,0.0,,
6 5,1.52379,0.0,,
```



and the fine-tuned model on a test set, and comparing the samples side by side. The test set should ideally include the full distribution of inputs that you might send to the model in a production use case. If manual evaluation is too time-consuming, consider using our Evals library to automate future evaluations.

# Iterating on data quality

If the results from a fine-tuning job are not as good as you expected, consider the following ways to adjust the training dataset:

Collect examples to target remaining issues

If the model still isn't good at certain aspects, add training examples that directly show the model how to do these aspects correctly

Scrutinize existing examples for issues

If your model has grammar, logic, or style issues, check if your data has any of the same issues. For instance, if the model now says "I will schedule this meeting for you" (when it shouldn't), see if existing examples teach the model to say it can do new things that it can't do

Consider the balance and diversity of data

If 60% of the assistant responses in the data says "I cannot answer this", but at inference time only 5% of responses should say that, you will likely get an overabundance of refusals

Make sure your training examples contain all of the information needed for the response

If we want the model to compliment a user based on their personal traits and a training example includes assistant compliments for traits not found in the preceding conversation, the model may learn to hallucinate information

Look at the agreement / consistency in the training examples

If multiple people created the training data, it's likely that model performance will be limited by the level of agreement / consistency between people. For instance, in a text extraction task, if people only agreed on 70% of extracted snippets, the model would likely not be able to do better than this

Make sure your all of your training examples are in the same format, as expected for inference

# Iterating on data quantity

Once you're satisfied with the quality and distribution of the examples, you can consider scaling up the number of training examples. This tends to help the model learn the task better, especially



increasing the training data size by:

Fine-tuning on your current dataset

Fine-tuning on half of your current dataset

Observing the quality gap between the two

In general, if you have to make a trade-off, a smaller amount of high-quality data is generally more effective than a larger amount of low-quality data.

# **Iterating on hyperparameters**

We allow you to specify the following hyperparameters:

epochs

learning rate multiplier

batch size

We recommend initially training without specifying any of these, allowing us to pick a default for you based on dataset size, then adjusting if you observe the following:

If the model does not follow the training data as much as expected increase the number of epochs by 1 or 2

This is more common for tasks for which there is a single ideal completion (or a small set of ideal completions which are similar). Some examples include classification, entity extraction, or structured parsing. These are often tasks for which you can compute a final accuracy metric against a reference answer.

If the model becomes less diverse than expected decrease the number of epochs by 1 or 2

This is more common for tasks for which there are a wide range of possible good completions

If the model does not appear to be converging, increase the learning rate multiplier

You can set the hyperparameters as is shown below:

```
from openai import OpenAI
client = OpenAI()

client.fine_tuning.jobs.create(
```



```
8 "n_epochs":2
9 }
```

# Fine-tuning examples

Now that we have explored the basics of the fine-tuning API, let's look at going through the fine-tuning lifecycle for a few different use cases.

- > Style and tone
- > Structured output
- > Function calling

# Migration of legacy models

For users migrating from /v1/fine-tunes to the updated /v1/fine\_tuning/jobs API and newer models, the main difference you can expect is the updated API. The legacy prompt completion pair data format has been retained for the updated babbage-002 and davinci-002 models to ensure a smooth transition. The new models will support fine-tuning with 4k token context and have a knowledge cutoff of September 2021.

For most tasks, you should expect to get better performance from <a href="gpt-3.5-turbo">gpt-3.5-turbo</a> than from the GPT base models.

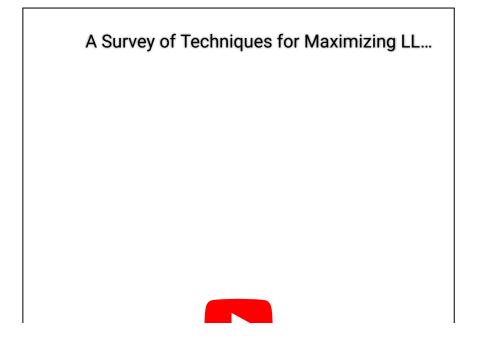
# **FAQ**

When should I use fine-tuning vs embeddings / retrieval augmented generation?



By default OpenAl's models are trained to be helpful generalist assistants. Fine-tuning can be used to make a model which is narrowly focused, and exhibits specific ingrained behavior patterns. Retrieval strategies can be used to make new information available to a model by providing it with relevant context before generating its response. Retrieval strategies are not an alternative to fine-tuning and can in fact be complementary to it.

You can explore the differences between these options further in our Developer Day talk:



#### Can I fine-tune GPT-4 or GPT-3.5-Turbo-16k?

GPT-4 fine-tuning is in experimental access and eligible developers can request access via the fine-tuning UI. Currently, [gpt-3.5-turbo-1106] supports up to 16K context examples.

# How do I know if my fine-tuned model is actually better than the base model?

We recommend generating samples from both the base model and the fine-tuned model on a test set of chat conversations, and comparing the samples side by side. For more comprehensive evaluations, consider using the OpenAl evals framework to create an eval specific to your use case.

# Can I continue fine-tuning a model that has already been fine-tuned?

Yes, you can pass the name of a fine-tuned model into the model parameter when creating a fine-tuning job. This will start a new fine-tuning job using the fine-tuned model as the starting



# How can I estimate the cost of fine-tuning a model?

Please refer to the estimate cost section above.

# Does the new fine-tuning endpoint still work with Weights & Biases for tracking metrics?

No, we do not currently support this integration but are working to enable it in the near future.

# How many fine-tuning jobs can I have running at once?

Please refer to our rate limit guide for the most up to date information on the limits.

#### How do rate limits work on fine-tuned models?

A fine-tuned model pulls from the same shared rate limit as the model it is based off of. For example, if you use half your TPM rate limit in a given time period with the standard gpt-3.5-turbo model, any model(s) you fine-tuned from gpt-3.5-turbo would only have the remaining half of the TPM rate limit accessible since the capacity is shared across all models of the same type.

Put another way, having fine-tuned models does not give you more capacity to use our models from a total throughput perspective.