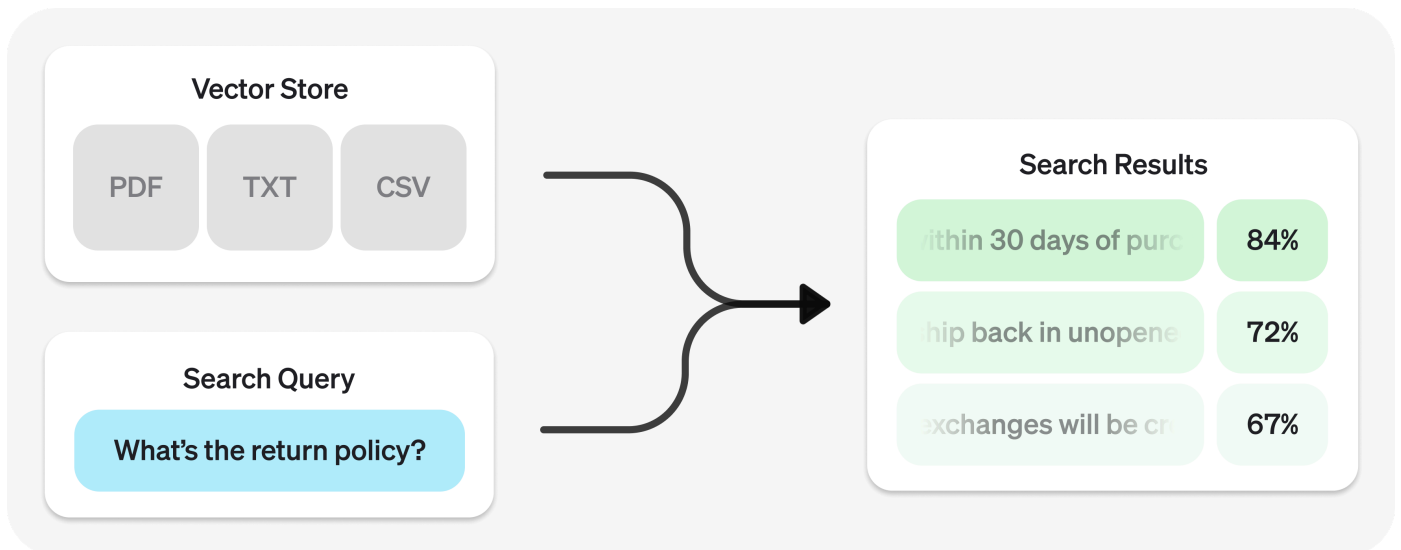


Retrieval

Search your data using semantic similarity.

The **Retrieval API** allows you to perform semantic search over your data, which is a technique that surfaces semantically similar results — even when they match few or no keywords. Retrieval is useful on its own, but is especially powerful when combined with our models to synthesize responses.



The Retrieval API is powered by vector stores, which serve as indices for your data. This guide will cover how to perform semantic search, and go into the details of vector stores.

Quickstart

- 1 Create vector store and upload files.

Create vector store with files

python

```
1 from openai import OpenAI
2 client = OpenAI()
3
4 vector_store = client.vector_stores.create(           # Create vector store
5     name="Support FAQ",
6 )
```

```

7
8 client.vector_stores.files.upload_and_poll(           # Upload file
9     vector_store_id=vector_store.id,
10    file=open("customer_policies.txt", "rb")
11 )

```

2 **Send search query** to get relevant results.

Search query

python ↕ 

```

1 user_query = "What is the return policy?"
2
3 results = client.vector_stores.search(
4     vector_store_id=vector_store.id,
5     query=user_query,
6 )

```

 To learn how to use the results with our models, check out the [synthesizing responses](#) section.

Semantic search

Semantic search is a technique that leverages [vector embeddings](#) to surface semantically relevant results. Importantly, this includes results with few or no shared keywords, which classical search techniques might miss.

For example, let's look at potential results for "When did we go to the moon?" :

TEXT	KEYWORD SIMILARITY	SEMANTIC SIMILARITY
The first lunar landing occurred in July of 1969.	0%	65%
The first man on the moon was Neil Armstrong.	27%	43%
When I ate the moon cake, it was delicious.	40%	28%

([Jaccard](#) used for keyword, [cosine with](#) `text-embedding-3-small` used for semantic.)

Notice how the most relevant result contains none of the words in the search query. This flexibility makes semantic search a very powerful technique for querying knowledge bases of any size.

Semantic search is powered by [vector stores](#), which we cover in detail later in the guide. This section will focus on the mechanics of semantic search.

Performing semantic search

You can query a vector store using the `search` function and specifying a `query` in natural language. This will return a list of results, each with the relevant chunks, similarity scores, and file of origin.

Search query

python ↕



```
1 results = client.vector_stores.search(  
2     vector_store_id=vector_store.id,  
3     query="How many woodchucks are allowed per passenger?",  
4 )
```

Results



```
1 {  
2     "object": "vector_store.search_results.page",  
3     "search_query": "How many woodchucks are allowed per passenger?",  
4     "data": [  
5         {  
6             "file_id": "file-12345",  
7             "filename": "woodchuck_policy.txt",  
8             "score": 0.85,  
9             "attributes": {  
10                "region": "North America",  
11                "author": "Wildlife Department"  
12            },  
13            "content": [  
14                {  
15                    "type": "text",  
16                    "text": "According to the latest regulations, each passenger is allow  
17                },  
18                {  
19                    "type": "text",  
20                    "text": "Ensure that the woodchucks are properly contained during tra  
21                }  
22            ]  
23        },  
24        {  
25            "file_id": "file-67890",  
26            "filename": "transport_guidelines.txt",  
27            "score": 0.75,  
28            "attributes": {  
29                "region": "North America",  
30                "author": "Transport Authority"  
31        },
```

```

32     "content": [
33         {
34             "type": "text",
35             "text": "Passengers must adhere to the guidelines set forth by the Tr
36         }
37     ]
38 }
39 ],
40 "has_more": false,
41 "next_page": null
42 }

```

A response will contain 10 results maximum by default, but you can set up to 50 using the `max_num_results` param.

Query rewriting

Certain query styles yield better results, so we've provided a setting to automatically rewrite your queries for optimal performance. Enable this feature by setting `rewrite_query=true` when performing a `search` .

The rewritten query will be available in the result's `search_query` field.

ORIGINAL	REWRITTEN
I'd like to know the height of the main office building.	primary office building height
What are the safety regulations for transporting hazardous materials?	safety regulations for hazardous materials
How do I file a complaint about a service issue?	service complaint filing process

Attribute filtering

Attribute filtering helps narrow down results by applying criteria, such as restricting searches to a specific date range. You can define and combine criteria in `attribute_filter` to target files based on their attributes before performing semantic search.

Use **comparison filters** to compare a specific `key` in a file's `attributes` with a given `value` , and **compound filters** to combine multiple filters using `and` and `or` .

Comparison filter



```

1 {
2   "type": "eq" | "ne" | "gt" | "gte" | "lt" | "lte" | "in" | "nin", // comparison
3   "key": "attributes_key", // attributes key
4   "value": "target_value" // value to compare against
5 }

```

Compound filter



```

1 {
2   "type": "and" | "or", // logical operators
3   "filters": [...]
4 }

```

Below are some example filters.

Region

Date range

Filenames

Exclude filenames

Complex

Filter for a region



```

1 {
2   "type": "eq",
3   "key": "region",
4   "value": "us"
5 }

```

Ranking

If you find that your file search results are not sufficiently relevant, you can adjust the `ranking_options` to improve the quality of responses. This includes specifying a `ranker`, such as `auto` or `default-2024-08-21`, and setting a `score_threshold` between 0.0 and 1.0. A higher `score_threshold` will limit the results to more relevant chunks, though it may exclude some potentially useful ones.

Vector stores

Vector stores are the containers that power semantic search for the Retrieval API and the [file search](#) tool. When you add a file to a vector store it will be automatically chunked, embedded, and indexed.


Vector stores contain `vector_store_file` objects, which are backed by a `file` object.

OBJECT TYPE	DESCRIPTION
<code>file</code>	Represents content uploaded through the Files API . Often used with vector stores, but also for fine-tuning and other use cases.
<code>vector_store</code>	Container for searchable files.
<code>vector_store.file</code>	Wrapper type specifically representing a file that has been chunked and embedded, and has been associated with a <code>vector_store</code> . Contains attributes map used for filtering.

Pricing

You will be charged based on the total storage used across all your vector stores, determined by the size of parsed chunks and their corresponding embeddings.

STORAGE	COST
Up to 1 GB (across all stores)	Free
Beyond 1 GB	\$0.10/GB/day

 See [expiration policies](#) for options to minimize costs.

Vector store operations

Create Retrieve Update Delete List

Create vector store

python 

```
1 client.vector_stores.create(  
2     name="Support FAQ",  
3     file_ids=["file_123"]  
4 )
```

Vector store file operations

Some operations, like `create` for `vector_store.file`, are asynchronous and may take time to complete — use our helper functions, like `create_and_poll` to block until it is. Otherwise, you

may check the status.

Create Upload Retrieve Update Delete List

Create vector store file

python ↕ 

```
1 client.vector_stores.files.create_and_poll(  
2     vector_store_id="vs_123",  
3     file_id="file_123"  
4 )
```

Batch operations

Create Retrieve Cancel List

Batch create operation

python ↕ 

```
1 client.vector_stores.file_batches.create_and_poll(  
2     vector_store_id="vs_123",  
3     file_ids=["file_123", "file_456"]  
4 )
```

Attributes

Each `vector_store.file` can have associated `attributes`, a dictionary of values that can be referenced when performing [semantic search](#) with [attribute filtering](#). The dictionary can have at most 16 keys, with a limit of 256 characters each.

Create vector store file with attributes

python ↕ 

```
1 client.vector_stores.files.create(  
2     vector_store_id="<vector_store_id>",  
3     file_id="file_123",  
4     attributes={  
5         "region": "US",  
6         "category": "Marketing",  
7         "date": 1672531200 # Jan 1, 2023  
8     }  
9 )
```

Expiration policies

You can set an expiration policy on `vector_store` objects with `expires_after`. Once a vector store expires, all associated `vector_store.file` objects will be deleted and you'll no longer be charged for them.

Set expiration policy for vector store

python ↕ 

```
1 client.vector_stores.update(  
2     vector_store_id="vs_123",  
3     expires_after={  
4         "anchor": "last_active_at",  
5         "days": 7  
6     }  
7 )
```

Limits

The maximum file size is 512 MB. Each file should contain no more than 5,000,000 tokens per file (computed automatically when you attach a file).

Chunking

By default, `max_chunk_size_tokens` is set to `800` and `chunk_overlap_tokens` is set to `400`, meaning every file is indexed by being split up into 800-token chunks, with 400-token overlap between consecutive chunks.

You can adjust this by setting `chunking_strategy` when adding files to the vector store. There are certain limitations to `chunking_strategy`:

`max_chunk_size_tokens` must be between 100 and 4096 inclusive.

`chunk_overlap_tokens` must be non-negative and should not exceed `max_chunk_size_tokens / 2`.

> Supported file types

Synthesizing responses

After performing a query you may want to synthesize a response based on the results. You can leverage our models to do so, by supplying the results and original query, to get back a grounded

response.

Perform search query to get results

python ↕



```
1 from openai import OpenAI
2
3 client = OpenAI()
4
5 user_query = "What is the return policy?"
6
7 results = client.vector_stores.search(
8     vector_store_id=vector_store.id,
9     query=user_query,
10 )
```

Synthesize a response based on results

python ↕



```
1 formatted_results = format_results(results.data)
2
3 '\n'.join('\n'.join(c.text) for c in result.content for result in results.data)
4
5 completion = client.chat.completions.create(
6     model="gpt-4.1",
7     messages=[
8         {
9             "role": "developer",
10            "content": "Produce a concise answer to the query based on the prov
11        },
12        {
13            "role": "user",
14            "content": f"Sources: {formatted_results}\n\nQuery: '{user_query}'"
15        }
16    ],
17 )
18
19 print(completion.choices[0].message.content)
```

"Our return policy allows returns within 30 days of purchase."



This uses a sample `format_results` function, which could be implemented like so:

Sample result formatting function

python ↕



```
1 def format_results(results):
2     formatted_results = ''
3     for result in results.data:
4         formatted_result = f"<result file_id='{result.file_id}' file_name='{result.file_name}'>"
5         for part in result.content:
6             formatted_result += f"<content>{part.text}</content>"
7         formatted_results += formatted_result + "</result>"
8     return f"<sources>{formatted_results}</sources>"
```
