

[Docs Home](#) / [MongoDB Atlas](#) / [Atlas Vector Search](#)

Run Vector Search Queries

On this page

Definition

[ANN Search](#)[ENN Search](#)

Syntax

Fields

Behavior

[Atlas Vector Search Index](#)[Atlas Vector Search Score](#)[Atlas Vector Search Pre-Filter](#)[Considerations](#)[Limitations](#)[Supported Clients](#)[Parallel Query Execution Across Segments](#)

Examples

ANN Examples

[ENN Example](#)

Atlas Vector Search queries take the form of an **aggregation pipeline stage**. For the `$vectorSearch` queries, Atlas Vector Search returns the results of your semantic search.

NOTE

Definition

The `$vectorSearch` stage performs an ANN or ENN search on a vector in the specified field.

ANN Search

For ANN search, Atlas Vector Search finds vector embeddings in your data that are closest to the vector embedding in your query based on their proximity in multi-dimensional space and based on the number of neighbors that it considers. It uses the **Hierarchical Navigable Small Worlds** algorithm and finds the vector embeddings most similar to the vector embedding in your query without scanning every vector. Therefore, ANN search is ideal for querying large datasets without significant filtering.

ENN Search

For ENN search, Atlas Vector Search exhaustively searches all the indexed vector embeddings by calculating the distance between all the embeddings and finds the exact nearest neighbor for the vector embedding in your query. This is computationally intensive and might negatively impact query latency. Therefore, we recommend ENN searches for the following use-cases:

- You want to determine the recall and accuracy of your ANN query using the ideal, exact results for the ENN query.
- You want to query less than 10000 documents without having to tune the number of nearest neighbors to consider.
- You want to include selective pre-filters in your query against collections where less than 5% of your data meets the given pre-filter.

Syntax

The field that you want to search must be indexed as Atlas Vector Search **vector** type inside a **vectorSearch** index type.

`$vectorSearch`

A `$vectorSearch` pipeline has the following prototype form:

```
{
  "$vectorSearch": {
    "exact": true | false,
    "filter": {<filter-specific>},
    "index": "<index-name>",
    "limit": <number-of-results>,
    "numCandidates": <number-of-
    "path": "<field-to-search>",
    "queryVector": [<array-of-nu
  }
}
```

Fields

The `$vectorSearch` stage takes a document with the following fields:

Field	Type	Necessity	Description
<code>exact</code>	boolean	Optional	<p>This is required if <code>numCandidates</code> is omitted.</p> <p>Flag that specifies whether to run ENN or ANN search. Value can be one of the following:</p> <ul style="list-style-type: none"><code>false</code> - to run ANN search<code>true</code> - to run ENN search <p>If omitted, defaults to <code>false</code>.</p> <p>Atlas Vector Search supports ANN search on clusters running MongoDB v6.0.11, v7.0.2, or later and ENN search on clusters running MongoDB v6.0.16, v7.0.10, v7.3.2, or later.</p>

Field	Type	Necessity	Description
<code>filter</code>	document	Optional	Any MQL match expression that compares an indexed field with a boolean, date, objectId, number (not decimals), string, or UUID to use as a pre-filter. To learn which query and aggregation pipeline operators Atlas Vector Search supports in your filter, see Atlas Vector Search Pre-Filter.
<code>index</code>	string	Required	<p>Name of the Atlas Vector Search index to use.</p> <p>Atlas Vector Search doesn't return results if you misspell the index name or if the specified index doesn't already exist on the cluster.</p>
<code>limit</code>	number	Required	Number (of type <code>int</code> only) of documents to return in the results. This value can't exceed the value of <code>numCandidates</code> if you specify <code>numCandidates</code> .

Field	Type	Necessity	Description
numCandidates	number	Optional	<p>This field is required if <code>exact</code> is <code>false</code> or omitted.</p> <p>Number of nearest neighbors to use during the search. Value must be less than or equal to (\leq) <code>10000</code>. You can't specify a number less than the number of documents to return (<code>limit</code>).</p> <p>We recommend that you specify a number higher than the number of documents to return (<code>limit</code>) to increase accuracy although this might impact latency. For example, we recommend a ratio of ten to twenty nearest neighbors for a limit of only one document. This overrequest pattern is the recommended way to trade off latency and recall in your ANN searches, and we recommend tuning this on your specific dataset.</p>
path	string	Required	Indexed vector type field to search.

Field	Type	Necessity	Description
queryVector	array of numbers	Required	<p>Array of numbers of the BSON <code>double</code>, BSON <code>BinData</code> <code>vector</code> subtype <code>float32</code>, or BSON <code>BinData</code> <code>vector</code> subtype <code>int1</code> or <code>int8</code> type that represent the query vector. The number type must match the indexed field value type. Otherwise, Atlas Vector Search doesn't return any results or errors.</p> <p>To learn more about generating BSON <code>BinData</code> <code>vector</code> subtype <code>float32</code> or <code>vector</code> subtype <code>int1</code> or <code>int8</code> vectors for your query, see How to Ingest Pre-Quantized Vectors.</p> <p>The array size must match the number of vector <code>dimensions</code> specified in the index definition for the field.</p> <p>You must embed your query with the same model that you used to embed the data.</p>

Behavior

`$vectorSearch` must be the first stage of any pipeline where it appears.

Atlas Vector Search Index

You must index the fields to search using the `$vectorSearch` stage inside a `vectorSearch` type index definition. You can index the following

types of fields in an Atlas Vector Search

vectorSearch type index definition:

- Fields that contain vector embeddings as **vector** type.
- Fields that contain boolean, date, objectId, numeric, string, and UUID values as **filter** type to enable vector search on pre-filtered data.

To learn more about these Atlas Vector Search field types, see [How to Index Fields for Vector Search](#).

Atlas Vector Search Score

Atlas Vector Search assigns a score, in a fixed range from 0 to 1 (where 0 indicates low similarity and 1 indicates high similarity), to every document that it returns. For `cosine` and `dotProduct` similarities, Atlas Vector Search normalizes the score to ensure that the score is not negative.

For `cosine`, Atlas Vector Search uses the following algorithm to normalize the score:

$$\text{score} = (1 + \text{cosine}(v1, v2)) / 2$$

For `dotProduct`, Atlas Vector Search uses the following algorithm to normalize the score:

$$\text{score} = (1 + \text{dotProduct}(v1, v2)) / 2$$

These algorithms show that Atlas Vector Search normalizes the score by considering the similarity score of the document vector ($v1$) and the query vector ($v2$), which has the range $[-1, 1]$. Atlas Vector Search adds 1 to the similarity score to normalize the score to a range $[0, 2]$ and then divides by 2 to ensure a value between 0 and 1.

For `euclidean` similarity, Atlas Vector Search uses the following algorithm to normalize the score to ensure a value between 0 and 1:

$$\text{score} = 1 / (1 + \text{euclidean}(v1, v2))$$

The preceding algorithm shows that Atlas Vector Search normalizes the score by calculating the euclidean distance, which is the distance between the document vector (v_1) and the query vector (v_2), which has the range $[0, \infty]$. Atlas Vector Search then transforms the distance to a similarity score by adding 1 to the distance and then divides 1 by the similarity score to ensure a value between 0 and 1 .

The score assigned to a returned document is part of the document's metadata. To include each returned document's score along with the result set, use a `$project` stage in your aggregation pipeline.

To retrieve the score of your Atlas Vector Search query results, use `vectorSearchScore` as the value in the `$meta` expression. That is, after the `$vectorSearch` stage, in the `$project` stage, the `score` field takes the `$meta` expression. The expression requires the `vectorSearchScore` value to return the score of documents for the vector search.

EXAMPLE

```

1  db.<collection>.aggregate([
2    {
3      "$vectorSearch": {
4        <query-syntax>
5      }
6    },
7    {
8      "$project": {
9        "<field-to-include>": 1,
10       "<field-to-exclude>": 0,
11       "score": { "$meta": "vectorSearchScore" }
12     }
13   }
14 ])
```

NOTE

Pre-filtering your data doesn't affect the score that Atlas Vector Search returns using `vectorSearchScore` for `vectorSearch` queries.

Atlas Vector Search Pre-Filter

The `$vectorSearch` `filter` option matches only BSON boolean, date, objectId, numeric, string, and UUID values. You **must** index the fields that you want to filter your data by as the `filter` type in a `vectorSearch` type index definition. Filtering your data is useful to narrow the scope of your semantic search and ensure that not all vectors are considered for comparison.

Atlas Vector Search supports the `$vectorSearch` `filter` option for the following MQL match expressions:

- `$gt`
- `$lt`
- `$gte`
- `$lte`
- `$eq`
- `$ne`
- `$in`
- `$nin`
- `$nor`
- `$not`
- `$and`
- `$or`

The `$vectorSearch` `filter` option supports only the following aggregation pipeline operators:

- `$and`
- `$or`

NOTE

The `$vectorSearch` `filter` option doesn't support other query operators, aggregation pipeline operators, or Atlas Search operators.

Considerations

Atlas Vector Search supports the short form of `$eq`. In the short form, you don't need to specify `$eq` in the query.

EXAMPLE

For example, consider the following filter with `$eq`:

```
"filter": { "_id": { "$eq": ObjectId("5a9427648" )
```

You can run the preceding query using the short form of `$eq` in the following way:

```
"filter": { "_id": ObjectId("5a9427648" )
```

You can also specify an array of filters in a single query by using the `$and` operator.

EXAMPLE

For example, consider the following pre-filter for documents with a `genres` field equal to `Action` and a `year` field with the value `1999`, `2000`, or `2001`:

```
"filter": {
  "$and": [
    { "genres": "Action" },
    { "year": { "$in": [ 1999, 2000, 2001 ] } }
  ]
}
```

Limitations

`$vectorSearch` is supported only on Atlas clusters running the following MongoDB versions:

- v6.0.11
- v7.0.2 and later (including RCs).

`$vectorSearch` can't be used in **view definition** and the following pipeline stages:

- `$lookup` sub-pipeline ★
- `$unionWith` sub-pipeline ★
- `$facet` pipeline stage

★ You can pass the results of `$vectorSearch` to this stage.

Supported Clients

You can run `$vectorSearch` queries by using the Atlas UI, `mongosh`, and any MongoDB driver.

You can also use Atlas Vector Search with local Atlas deployments that you create with the Atlas CLI. To learn more, see [Create a Local Atlas Deployment](#).

Parallel Query Execution Across Segments

We recommend dedicated **Search Nodes** to isolate vector search query processing. You might see improved query performance on the dedicated Search Nodes. Note that the high-CPU systems might provide more performance improvement. When Atlas Vector Search runs on search nodes, Atlas Vector Search parallelizes query execution across segments of data.

Parallelization of query processing improves the response time in many cases, such as queries on large datasets. Using intra-query parallelism during Atlas Vector Search query processing utilizes more resources, but improves latency for each individual query.

NOTE

Atlas Vector Search doesn't guarantee that each query will run concurrently. For example, when too many concurrent queries are queued, Atlas Vector Search might fall back to single-threaded execution.

You might see inconsistent results for the same successive queries. To mitigate this, increase the value of `numCandidates` in your `$vectorSearch` queries.

Examples

The following queries search the sample `sample_mflix.embedded_movies` collection using the `$vectorSearch` stage. The queries search the `plot_embedding` field, which contains embeddings created using OpenAI's `text-embedding-ada-002` embeddings model.

If you added the **sample** collection to your Atlas cluster and created the **sample** indexes for the collection, you can run the following queries against the collection.

NOTE

Pasting the `queryVector` from the sample code into your terminal might take a while depending on your machine.

► Use the **Select your language** drop-down menu to set the language of the examples in this page.


ANN Examples

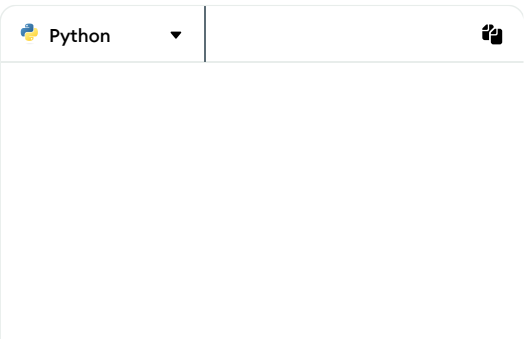
Basic Example**Filter Example**

The following query uses the `$vectorSearch` stage to search the `plot_embedding` field using vector embeddings for the string *time travel*. It considers up to `150` nearest neighbors, and returns `10` documents in the results. The query also specifies a `$project` stage to do the following:

- Exclude the `_id` field and include only the `plot` and `title` fields in the results.
- Add a field named `score` that shows the **vector search score** for each document in the results.

TIP

Work with a runnable version of this example as a [Python notebook](#). 



```

1  import pymongo
2
3  # connect to your Atlas cluster
4  client = pymongo.MongoClient("conne
5
6  # define pipeline
7  pipeline = [
8      {
9          '$vectorSearch': {
10             'index': 'vector_index',
11             'path': 'plot_embedding',
12             'queryVector': [-0.0016261312,
13             'numCandidates': 150,
14             'limit': 10
15         }
16     }, {
17         '$project': {
18             '_id': 0,
19             'plot': 1,
20             'title': 1,
21             'score': {
22                 '$meta': 'vectorSearchScore'
23             }
24         }
25     }
26 ]
27
28 # run pipeline
29 result = client["sample_mflix"]["eml
30
31 # print results
32 for i in result:
33     print(i)
34

```

^ HIDE OUTPUT

```

1  {'plot': 'A reporter, learning of t
2  {'plot': 'At the age of 21, Tim disc
3  {'plot': 'Hoping to alter the events
4  {'plot': 'After using his mother's r
5  {'plot': 'An officer for a security
6  {'plot': 'A time-travel experiment
7  {'plot': 'Agent J travels in time to
8  {'plot': 'Bound by a shared destiny,
9  {'plot': 'With the help of his uncl
10 {'plot': 'A dimension-traveling wiza

```

ENN Example

The following query uses the `$vectorSearch` stage to search the `plot_embedding` field using vector embeddings for the string *world war*. It requests exact matches and limits the results to 10 documents only. The query also specifies a `$project` stage to do the following:

- Exclude the `_id` field and include only the `plot`, `title`, and `year` fields in the results.
- Add a field named `score` that shows the vector search score of the documents in the results.

TIP


Work with a runnable version of this example as a [Python notebook](#).[↗]

 Python ▼



```
1  import pymongo
2
3  # connect to your Atlas cluster
4  client = pymongo.MongoClient("<conn
5
6  # define pipeline
7  pipeline = [
8      {
9          '$vectorSearch': {
10             'index': 'vector_index',
11             'path': 'plot_embedding',
12             'queryVector': [-0.006954097, -
13             'exact': True,
14             'limit': 10
15         }
16     }, {
17         '$project': {
18             '_id': 0,
19             'plot': 1,
20             'title': 1,
21             'score': {
22                 '$meta': 'vectorSearchScore'
23             }
24         }
25     }
26 ]
27
28 # run pipeline
29 result = client["sample_mflix"]["eml
30
31 # print results
32 for i in result:
33     print(i)
34
```

 HIDE OUTPUT

 MongoDB Docs

MongoDB Atlas

► Get Started

► Create & Connect to Clusters

► Configure Security Features

► Configure UI Access

► Migrate or Import Data

► Interact with Data

► Query Federated Data

► Atlas Search

▼ Atlas Vector Search

Quick Start

Create Embeddings

Create and Manage Indexes

Create and Run Queries

Vector Quantization

Retrieval-Augmented Generation (RAG)

Review Deployment Options

► Tutorials

► AI Integrations

Evaluate Results

Improve Performance

```
1  {
2    'plot': 'It is the dawn of World War I',
3    'title': 'Red Dawn',
4    'score': 0.7700583338737488
5  }
6  {
7    'plot': 'A dramatization of the World War I',
8    'title': 'Sands of Iwo Jima',
9    'score': 0.7581185102462769
10 }
11 {
12   'plot': 'Great Patriotic War, early years',
13   'title': 'White Tiger',
14   'score': 0.750884473323822
15 }
16 {
17   'plot': 'As World War Two rages on, a P-51 Dragon Fighter',
18   'title': 'P-51 Dragon Fighter',
19   'score': 0.749922513961792
20 }
21 {
22   'plot': 'A private in the latter part of World War II',
23   'title': 'When Trumpets Fade',
24   'score': 0.7498313188552856
25 }
26 {
27   'plot': 'Post World War III future',
28   'score': 0.7497193217277527
29 }
30 {
31   'plot': 'It is post-World War III. A robot named Jox',
32   'title': 'Robot Jox',
33   'score': 0.7495121955871582
34 }
35 {
36   'plot': 'During World War II, an American soldier',
37   'title': 'The Enemy Below',
38   'score': 0.746050238609314
39 }
40 {
41   'plot': 'Four American soldiers and a Japanese soldier',
42   'title': 'Saints and Soldiers',
43   'score': 0.743497371673584
44 }
45 }
```


Troubleshooting

Changelog

▶ Atlas Stream Processing

▶ Backup, Restore, and Archive

▶ Resource Tags

▶ Manage Clusters

▶ Monitor Clusters

▶ Related Services

▶ Manage Billing

▶ Programmatic Access

▶ Atlas CLI

upcomi... ▼

About

Careers

Legal

Security Information

Connect with Us

Support

Contact Us

Atlas Status

Manage Cookies

Deployment Options

MongoDB Atlas

Community Edition

© 2024 MongoDB, Inc.

English

Investor Rela

GitHub

Trust Center

Customer Po

Customer Su

Enterprise Ac