

# ベクトルDB管理APIチュートリアル

## はじめに

本書は、本サービスが提供するベクトルDB管理APIのチュートリアルを記載します。

また、本書の対象読者は以下の通りです。

- ・本サービスと連携したシステムや製品開発を行う開発者
- ・管理ポータルを使用せず、ベクトル検索DBのリソースを操作したいユーザ

なお、本書ではPythonのサンプルコードを記載しておりますが、本サービスのAPIはREST形式のため、他の言語からでもご利用いただけます。

## ベクトルDB管理API チュートリアルの流れ

ベクトルDB管理APIでは、ベクトル検索DBに対して、インデックスの作成削除や、文書の登録削除を行う機能を有します。本チュートリアルでは、これらの一連の操作方法を記載します。


本チュートリアルの流れは以下です。

1. インデックスを作成する
2. (応用) requestsライブラリを利用したPythonプログラムで文書登録をする
3. (応用) LangChainライブラリを利用したPythonプログラムでチャンク文書を登録する
4. 登録した文書に対してチャンク検索を行う
5. 登録した文書に対して検索対話を行う
6. 登録した文書情報を取得する
7. 登録済みインデックス情報を取得する
8. インデックス使用量 (テナント合計・インデックス別) を取得する
9. 登録済みインデックス情報の変更を行う
10. 登録した文書情報を削除する
11. 作成したインデックスを削除する

## インデックスを作成する

コマンド実行でベクトル検索DB上に新規にインデックスを作成することが出来ます。

作成されたインデックスにはグループやユーザが紐づいていないので、引き続きコマンドからインデックスの操作を行いたい際は、作成後に管理ポータルから作成したインデックスにAPI USER、もしくはAPI USERが所属するグループの紐づけを行ってください。

-  上記のAPI USERとは事前にシステムに登録されているシステムユーザを指します。詳細は「管理ポータル操作ガイド (ユーザ登録編)」をご参照ください。

- ・インデックスにグループやユーザを紐づける方法は、「管理ポータル操作ガイド（インデックス・文書管理編）」をご参照ください。

## APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。API Keyは管理ポータルの契約情報から確認できます。

```
1 curl -X POST https://api.genai-api.nec-cloud.com/genai-search-api/index/createIndex ¥  
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>" ¥  
3 --data "<リクエストパラメータ>"
```

リクエストパラメータは必須のパラメータと任意のパラメータがあります。

パラメータ名	必須	説明
vectorIndex	○	作成するインデックス名を指定します。インデックス名には英数字(全角/半角)と日本語(全角/半角)とハイフン(-とー)のみ使用可能です(最初と最後にハイフンを使用することは不可)。また文字数は2~64文字の範囲で指定が可能です。加えて、すでに作成済みのインデックス名を指定することはできません。
embeddingModel	-	文書登録される際のベクトル化モデルをインデックスに紐づけて登録するパラメータで、管理ポータルの利用可能な埋め込みモデル一覧のモデルIDから確認できます。
comment	-	インデックスに対する説明文などをインデックスに紐づけて登録できるパラメータです。登録した内容はインデックス一覧画面から確認できます。

以下はリクエストパラメータの例です。

```
1 {  
2   "vectorIndex": "test-index",  
3   "embeddingModel": "multilingual-e5-large",  
4   "comment": "テスト用インデックス"  
5 }
```

インデックスの作成に成功するとJson形式で作成されたインデックスの名前が返却されます。

```
1 {"vectorIndex": "test-index"}
```

## (応用) requestsライブラリを利用したPythonプログラムで文書登録をする

文書登録には以下で紹介するpythonプログラムを利用することができます。

この章では、Base64形式にエンコードしたファイルをrequestsライブラリでAPIに送信して文書登録を行うPythonプログラム(以降、文書登録プログラム)を紹介します。

### **i** 文書登録時の注意事項

- ・ 契約したインデックスサイズを超過すると文書登録ができなくなりますので、事前にインデックスサイズをご確認ください。
- ・ 契約インデックスサイズの確認方法については、「管理ポータル操作ガイド（利用状況確認編）」をご参照ください。
- ・ 登録先インデックスについて、事前にAPI USERが所属しているグループ、もしくはAPI USERがインデックスに紐づいていることをご確認ください。
- ・ 拡張子が .txt や .csv のファイルを登録する際、内部で文字コードの特定ができず、登録エラーとなる場合があります。  
このような場合は、文字コードが判別しやすい形式（例：UTF-8 BOM付き）でファイルを保存し直したうえで、再度登録をお試しください。
- ・ ファイル内に含まれる「文字情報」のみを読み込み・解析対象とし、埋め込み画像・スキャン画像・図表画像内の文字情報を直接認識することはできません。

## requestsライブラリのインストール

文書登録プログラムでは、`requests` ライブラリなどが必要ですので、以下のコマンドでインストールを行ってください。

インストールコマンドの例

```
1 | pip install requests tzdata
```

## 文書登録プログラムのダウンロード

以下からrequestsライブラリを使用して、インデックスへの文書登録を行うPythonファイルを参照できるため、ローカルのPython実行環境に同じ記載のプログラムを配置してください。

左メニューから該当のプログラムを参照できます。

その他 > RAG文書登録プログラム

## 文書登録プログラムの実行 - パラメータ詳細

配置したPythonプログラムを実行することで指定したインデックスへの文書登録ができます。

実行コマンドには以下のようなパラメータがあり、必須のパラメータを指定して実行して下さい。また任意のパラメータについては `--<任意のパラメータ名> 値` のように指定して下さい。

パラメータ名	必須	説明
api_url	○	<p>文書登録APIのURL。管理ポータルAPI情報のAPIベースURLを確認して、”APIベースURL” + ” /genai-search-api/document/addDocument” を指定して下さい。</p> <p>また、非同期実行の場合は”APIベースURL” + ” /genai-search-api/document/addDocumentAsync” を指定して下さい（チャンク文書登録プログラムは対応していません）。</p>
auth_token	○	API認証に使用するトークン。管理ポータルAPI情報のAPIキーを指定して下さい。
directory_or_file_path	○	<p>処理対象のファイルまたはフォルダのパス</p> <ul style="list-style-type: none"> <li>• ファイルの場合：指定されたファイルを処理対象として読み込みます。</li> <li>• フォルダの場合：指定したフォルダ配下のファイルを再帰的に探索し、全てのファイルを処理対象として読み込みます。</li> </ul> <p>&lt;データの送信方式について&gt;</p> <ul style="list-style-type: none"> <li>• RAG文書登録プログラム <ul style="list-style-type: none"> <li>◦ 読み込んだファイルは base64エンコードされ、文書登録APIに送信されます。</li> </ul> </li> <li>• チャンク文書登録プログラム <ul style="list-style-type: none"> <li>◦ 読み込んだファイルは自動的にチャンク分割され、文書登録APIに送信されます。</li> <li>◦ 送信されるチャンクは以下のようなJson形式で構成されています： <ul style="list-style-type: none"> <li>▪ page_content : チャンクの中身（テキスト内容）。</li> </ul> </li> </ul> </li> </ul>

		<ul style="list-style-type: none"> <li>▪ <code>metadata</code> : チャンクに関連するメタデータ情報。</li> </ul> <p>&lt;APIに送信されるチャンクデータの例&gt;</p> <pre>{ "page_content": "ドキュメントの内容 1",   "metadata": { "source": "/aa/bb/cc.ppt",                 "page": 0 }}</pre>
<code>vector_index</code>	○	文書を登録するインデックス名
<code>url</code>	-	登録ファイルのURL。検索対話時にメタデータとして参照できるようになります。 <code>file_or_directory_path</code> でフォルダパスを指定時は、 <code>file_or_directory_path</code> の値をベースURLとして、それぞれのファイルパスに階層ごとに別々のパスを付与して登録を行います。
<code>overwrite</code>	-	上書きフラグ (Trueの場合、既存のファイルを上書きをします。)  デフォルトは上書きされます。
<code>custom_metadata</code>	-	任意のメタデータ (キーと値のペア) を指定できます。Json形式で指定し、検索対話時にメタデータとして参照できるようになります。
<code>kwargs</code>	-	文書登録時の追加オプション指定。 文書のチャンク化を実施する際のオプションを指定できます。  <code>split_chunk_size</code> : テキストチャンクのサイズ (文章を分割するときのサイズ) で、0~512の範囲で選択可能。指定なしの場合は500で動作します。  <code>split_overlap_size</code> : テキストチャンクのオーバーラップサイズ (各チャンクの境界が重なる部分のサイズ) で0~ <code>split_chunk_size</code> の指定値の範囲で選択可能。指定なしの場合は128で動作します。

以下、実行環境に応じたコマンド実行例になります。

## 実行コマンド(Linux)

```
1 python script_addDocuments.py <api_url> <auth_token> <directory_or_file_path> ¥
2 <vector_index> --url <base_url> --overwrite <overwrite> --custom_metadata <custom_data> ¥
3 --kwargs '{"split_chunk_size": "<chunk_size>", "split_overlap_size": "<overlap_size>"}
```

Windows環境のコマンドプロンプトでは仕様上、ダブルクォーテーションをエスケープする必要があるため、--kwargsオプション指定時は以下を参考にしてください。

## 実行コマンド(Windows - コマンドプロンプト)

```
1 python script_addDocuments.py <api_url> <auth_token> <directory_or_file_path> ^
2 <vector_index> --url <base_url> --overwrite <overwrite> --custom_metadata <custom_data> ^
3 --kwargs "{\"split_chunk_size\": \"¥<chunk_size>¥\", \"split_overlap_size\": \"¥<overlap_size>¥\"}"
```

登録に成功すると、以下のような実行結果が返却されます。

## 実行結果の例

```
1 Processing single file: D:¥test¥test.txt
2 Processed test.txt at 2024-07-23 11:54:43.013309+09:00: 200 null
```

## 非同期実行について

文書のサイズが大きい場合や多数のファイルを登録する際には、非同期実行を使用することを推奨します。非同期実行では、`/document/addDocumentAsync` エンドポイントを使用し、リクエストが完了するのを待たずに処理を進めることができます。これにより、大量の文書を効率的に登録することが可能です。ただし、1テナントごとに1日で非同期登録できるサイズには上限がありますので、ご注意ください。

## 登録に成功した場合の実行結果

文書の登録受付に成功すると、以下のような実行結果が返されます。

## 実行結果の例

```
1 Starting process for directory or file: D:¥test.txt
2 Processing single file: D:¥test.txt
3 Processed test.txt at 2024-10-16 10:29:23.522842+09:00: 200 {"sas_url":"https://~"}
```

非同期実行の場合には、上記の実行結果の例のようにURL ({"sas\_url":"https://~"}) が返されます。このURLを使用して処理の詳細を確認できます。URLにアクセスすると登録状況が記載されたJSONファイルを参照できます。

## 登録状況の確認

ブラウザなどで上記のURLにアクセスすることで、以下のようなJSONファイルを参照することができます。

```
1 {
```

```
2  "id": "74da4fb8-719f-4e08-884e-d64659d86171",
3  "timestamp": "2024-10-16T10:30:11Z",
4  "filename": "test.txt",
5  "status": "completed",
6  "details": {}
7  }
```

このJSONファイルには、登録されたファイルに関する情報が含まれています。

- **id**: 登録ファイルの一意のID。
- **timestamp**: 文書が登録された日時。
- **filename**: 登録対象のファイル名。
- **status**: 登録のステータス（受付が完了で「received」、登録処理が進行中で「processing」、完了で「completed」、失敗で「failed」）。
- **details**: 異常時にエラー情報が格納されます。

この情報を確認することで、文書登録の進行状況や結果を把握することができます。

## （応用）LangChainライブラリを利用したpythonプログラムでチャンク文書を登録する

この章では、LangChainというライブラリを使用して文書を読み込み、扱いやすいサイズに分割したチャンクを用いて文書登録を行うPythonプログラム(以降、チャンク文書登録プログラム)を紹介します。

今回紹介するのは、テキストファイルの読み取りと分割のみに対応したプログラムです。他の拡張子の文書の読み込み部分やチャンク分割部分は、実際のユースケースに合わせて適宜拡張してください。

### LangChainとは？

LangChainは、自然言語処理（NLP）や生成AIを利用したアプリケーションを簡単に作成するためのPythonライブラリです。このライブラリは以下のような特徴があります：

- AIモデルとの連携が簡単： LangChainはOpenAIのGPTなどのAIモデルを活用しやすいように設計されています。
- 大規模な文書进行处理する機能： 大量のテキストを効率的に処理し、AIでの応答生成や検索機能の実装が可能です。
- 「チャンク」と呼ばれる分割処理のサポート： 長い文書を自動的に分割し、AIモデルが扱いやすくする仕組みを提供しています。

---

### チャンクとは？

チャンク（chunk）とは、長い文章や文書を小さな部分に分割したものです。

通常、アプリケーションのAPI内部で文書は自動的に分割されますが、ユーザが分割のルールや細かい設定を自由にカスタマイズしたい場合、この機能を利用することで、分割方法を支配的にコントロールできます。以下のような場合に役立ちます：

- **特定の分割ルールを指定したい場合：**
  - 文の意味が途切れないように段落ごとに分割したい。
  - 固定の文字数で分割したい。
- **分割後の内容を事前に確認したい場合：**
  - 分割後のチャンクが適切かを確認し、必要に応じて微調整したい。

---

## 必要なライブラリのインストール

以下で紹介するチャンク文書登録プログラムでは、requestsライブラリのほかに、LangChain関連のライブラリが必要です。以下のコマンドを実行してインストールを行ってください。

```
1 | pip install langchain langchain-community requests chardet tiktoken tzdata transformers
```

## チャンク文書登録プログラムのダウンロード

以下からrequestsライブラリを使用して、インデックスへの文書登録を行うpythonファイルを参照できるため、ローカルのPython実行環境に同じ記載のプログラムを配置してください。また、本プログラムはテキストファイルの読み込みに対応しております。別の拡張子のファイルを読み取って実行する必要がある場合は、お手数ですが、以下のプログラムを参考に改修を行っていただきますようお願いいたします。

プログラムは左メニュー（その他 > チャンク文書登録プログラム）から参照できます。

また、実行時に関連ファイルが必要になるため、実行前に以下の作業を実施してください。

1. 「intfloat--multilingual-e5-large-tokenizer」という名前のフォルダを作成し、チャンク文書登録プログラムと同じ階層に配置してください。
2. [intfloat/multilingual-e5-large at main](#) から以下4つのファイルをダウンロードして、上記のフォルダに格納してください。

- config.json
- special\_tokens\_map.json
- tokenizer\_config.json
- tokenizer.json

## チャンク文書登録プログラムの実行

配置したpythonプログラムを実行することで指定したインデックスへのチャンク文書登録ができます。

実行コマンドは以下で、パラメータについてはRAG文書登録プログラムと同じになりますので、詳細は前章の文書登録プログラムの実行 - パラメータ詳細をご覧ください。

### 実行コマンド(Linux)

```
1 python script_addChunks.py <api_url> <auth_token> <directory_or_file_path> ¥
2 <vector_index> --url <base_url> --overwrite <overwrite> --custom_metadata <custom_data> ¥
3 --kwargs '{"split_chunk_size": "<chunk_size>", "split_overlap_size": "<overlap_size>"}'
```

Windows環境のコマンドプロンプトでは仕様上、ダブルクォーテーションをエスケープする必要がありますため、--kwargsオプション指定時は以下を参考にしてください。

### 実行コマンド(Windows - コマンドプロンプト)

```
1 python script_addChunks.py <api_url> <auth_token> <directory_or_file_path> ^
2 <vector_index> --url <base_url> --overwrite <overwrite> --custom_metadata <custom_data> ^
3 --kwargs "{¥"split_chunk_size¥": ¥"<chunk_size>¥", ¥"split_overlap_size¥": ¥"<overlap_size>¥"}"
```

登録に成功すると以下のような実行結果が返却されます。

### 実行結果の例

```
1 Processing single file: D:¥test¥test.txt
2 Processed test.txt at 2024-07-23 11:54:43.013309+09:00: 200 null
```

#### 実行時に表示される警告メッセージについて

プログラムの実行時に、以下のような警告メッセージが表示される場合がありますが、問題なく動作するため、このメッセージは無視していただいて構いません(学習や推論用のフレームワークが見つからないことを示す警告メッセージです。)

```
1 None of PyTorch, TensorFlow >= 2.0, or Flax have been found.
2 Models won't be available and only tokenizers, configuration and file/data utilities can be used.
```

## 登録した文書に対してチャンク検索を行う

以下のコマンド実行で、入力されたキーワードをもとに検索し、キーワードに関連するチャンク(文書の一部)を取得することができます。

また、文書登録時に付与したメタ情報(メタデータ)のKey/Value条件を指定して、検索対象を絞り込むことができます。

例えば、ファイル名やページ番号などのメタ情報を条件にし、条件に一致するチャンクのみを検索対象とすることが可能です。

## 対象となるメタデータ項目

メタデータフィルタで指定できる対象は以下です。

- `source` (例: ファイル名)
- `page` (例: ページ番号)
- `url` (例: 参照URL)
- ユーザが任意に指定する追加メタデータ

検索対象は、パラメータで指定したインデックス、かつAPI USERが所属するグループ、もしくはAPI USERが紐づいてるものが対象になります。

**i** メタデータフィルタは、登録時と同じ型・表記でないとは一致しません。

また、追加メタデータで絞込を行う場合は、日付のような値 (例: "2025-12-12") は入力と型解釈がずれてヒットしない場合があります。

そのため、必要に応じて `str_` や `date_` などの接頭辞を付け、型解釈の揺れを防いでください。

## APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。API Keyは管理ポータル契約情報から確認できます。

```
1 curl -X POST https://api.genai-api.nec-cloud.com/genai-search-api/document/searchRelatedChunks/ ¥  
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>" ¥  
3 --data "<リクエストパラメータ>"
```

リクエストパラメータは必須のパラメータと任意のパラメータがあります。

リクエストパラメータの各項目の詳細は以下のとおりです。

パラメータ名	型	デフォルト	必須	説明	備考
searchWord	string	-	○	検索する文字列	
vectorIndex	string	-	△	検索対象となるベクトルストアのインデックス名を指定する。	vectorIndex、vectorIndexesのどちらかを指定する。

vectorIndexes	array[string]	-	△	検索対象となる複数のベクトルストアのインデックス名を指定する場合は配列で指定する。	vectorIndex、vectorIndexesのどちらかを指定する。
searchOption	object	{"searchType": "hybrid"}	-	ベクトル検索のオプションをJson形式で必要に応じて複数記載する	
searchOption.searchType	string	hybrid	-	検索手法  similarity: 類似度で検索する方式  hybrid: 類似度検索とキーワード検索を組み合わせで検索する方式	
searchOption.topK	int	4	-	返却するチャンク数の最大数	最小値:1 最大値:5000※1
metadataFilters	array[object]	-	-	登録文書のメタ情報のKey/Valueを指定して検索対象チャンクを絞り込むための条件リスト (AND 条件)。	省略時はメタ情報による絞り込みは行わない。
metadataFilters[].key	string	-	○	マッチング対象とするメタデータのキー名。文書登録時の metadata に格納されている項目名を指定する。	キーが複数階層で登録されている場合は以下のように指定する。  例: 登録時 customMetadata: { "foo": { "bar": "x" } }  指定フィルタ: key="customMetadata.foo.bar", value="x"
metadataFilters[].value	str, int, float, bool	-	○	比較対象となる値。文書登録時のメタデータに保存さ	現時点でサポートしている完全

			れている値と同じ型・内容を指定する。	一致ではstr, int, float, boolの型のみで、それ以外は指定してもヒットしない。
--	--	--	--------------------	--

**i** `metadataFilters` は省略可能です。省略時はメタデータによる絞り込みは行わず、指定したインデックスに紐づく全文書が対象となります。

また、`metadataFilters` を複数指定した場合、条件はANDで評価されます。

**⚠** ※1: `searchOption.topK` は1~5000の範囲で指定可能ですが、100以下を推奨します。101以上を指定する場合は非推奨となります。リクエスト投入数に関する制約を含め、詳細はサービス仕様書の「10.機能およびAPIの諸元」を参照してください。

以下はリクエストパラメータの例です。

```

1  {
2  "searchWord": "スナップショット",
3  "vectorIndexes": ["test-index", "test-index2"],
4  "searchOption": {
5    "searchType": "hybrid",
6    "topK": 6
7  }
8  "metadataFilters": [
9    {
10   "key": "source",
11   "value": "test.pdf"
12  }
13  ]
14 }

```

検索に成功するとJson形式で関連するチャンクが関連度の高い順位に返却されます。

レスポンスボディの各項目の詳細は以下のとおりです。

パラメータ名	型	必須	説明
<code>relatedChunks</code>	<code>array[object]</code>	○	<code>searchWord</code> で指定された文章に近いチャンクの情報を返却する
<code>relatedChunks[].score</code>	<code>float</code>	○	マッチしたチャンクの関連度スコア。
<code>relatedChunks[].content</code>	<code>string</code>	○	チャンク本文

relatedChunks[].metadata	object	○	チャンクのメタデータ
relatedChunks[].metadata.source	string	○	チャンク元のファイルパス
relatedChunks[].metadata.folder	string	○	インデックス名
relatedChunks[].metadata.insertDate	string	○	登録時間
relatedChunks[].metadata.url	string	○	ソース文書の参照URL(文書登録時に任意に指定できる)
relatedChunks[].metadata.page	int	-	ソース文書のページ番号
relatedChunks[].metadata.任意のメタデータ	Object	-	文書登録時に任意のメタデータを登録している場合、任意のメタデータを返却する。

以下はレスポンスの例です。

```

1  {
2    "relatedChunks": [
3      {
4        "score": 0.9326973,
5        "content": "インデックス合算の最大登録容量は100GBです。管理ポータルからの確認方法は…",
6        "metadata": {
7          "source": "test.pdf",
8          "folder": "test-index",
9          "insertDate": "2025-10-01T08:30:12Z",
10         "url": "https://portal.example.com/docs/rag/guide_v2.pdf#page=12",
11         "page": 12
12       }
13     },
14     {
15       "score": 0.917397,
16       "content": "Snapshot retention の最小・最大保持数は…",
17       "metadata": {
18         "source": "test.pdf",
19         "folder": "test-index2",
20         "insertDate": "2025-09-28T04:11:47Z",
21         "url": "",
22         "page": 10
23       }
24     }
25   ]
26 }

```

**i** 検索でヒットする文書数が指定した値より小さい場合、返却する数はtopKで指定した数より小さくなる場合があります。

## 登録した文書による検索対話を利用する

文書登録を行ったインデックスを選択して、検索対話を利用することができます。利用方法は、チャットUIから利用する場合と検索対話のAPIから利用する場合の2パターンがあります。

- チャットUIから検索対話を行いたい場合は「チャット画面利用ガイド」をご確認ください。
- APIから検索対話を行いたい場合は「検索対話チュートリアル」をご確認ください。

## 登録した文書情報を取得する

以下のコマンド実行でインデックスに登録済みの文書の情報を取得できます。

また、文書登録時に付与したメタ情報（メタデータ）のKey/Value条件を指定して、検索対象を絞り込むことができます。

例えば、ファイル名や任意の追加メタデータなどのメタ情報を条件にし、条件に一致する文書のみを検索対象とすることが可能です。

### 対象となるメタデータ項目

メタデータフィルタで指定できる対象は以下です。

- `source`（例：ファイル名）
- `url`（例：参照URL）
- ユーザが任意に指定する追加メタデータ

実行時は指定するインデックスについて、事前にAPI USERが所属しているグループ、もしくはAPI USERがインデックスに紐づいていることをご確認ください。

**i** メタデータフィルタは、登録時と同じ型・表記でないと一致しません。

また、追加メタデータで絞込を行う場合は、日付のような値（例：“2025-12-12”）は入力と型解釈がずれてヒットしない場合があります。

そのため、必要に応じて `str_` や `date_` などの接頭辞を付け、型解釈の揺れを防いでください。

## APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。API Keyは管理ポータル契約情報から確認できます。

```
1 curl -X POST https://api.genai-api.nec-cloud.com/genai-search-api/document/ListDocument ¥  
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>" ¥
```

リクエストパラメータは必須のパラメータと任意のパラメータがあります。

リクエストパラメータの各項目の詳細は以下のとおりです。

パラメータ名	型	デフォルト	必須	説明	備考
vectorIndex	string	-	○	インデックス名	
metadataFilters	array[object]	-	-	登録文書のメタ情報のKey/Valueを指定して検索対象チャンクを絞り込むための条件リスト (AND 条件)。	省略時はメタ情報による絞り込みは行わない。
metadataFilters[].key	string	-	○	マッチング対象とするメタデータのキー名。文書登録時のmetadataに格納されている項目名を指定する。	キーが複数階層で登録されている場合は以下のように指定する。 例：登録時 customMetadata: { "foo": { "bar": "x" } }  指定フィルタ： key="customMetadata.foo.bar", value="x"
metadataFilters[].value	str, int, float, bool	-	○	比較対象となる値。文書登録時のメタデータに保存されている値と	現時点でサポートしている完全一致ではstr, int, float, boolの型の

			同じ型・内容を指定する。	みで、それ以外は指定してもヒットしない。
--	--	--	--------------	----------------------

**i** `metadataFilters` は省略可能です。省略時はメタデータによる絞り込みは行わず、指定したインデックスに紐づく全文書が対象となります。

また、`metadataFilters` を複数指定した場合、条件はANDで評価されます。

以下はリクエストパラメータの例です。

```

1 {
2   "vectorIndex": "test-index",
3   "metadataFilters": [
4     {
5       "key": "source",
6       "value": "test.txt"
7     }
8   ]
9 }
10 }
11

```

レスポンスボディの各項目の詳細は以下のとおりです。

パラメータ名	型	デフォルト	必須	説明	備考
<code>totalCount</code>	<code>int</code>	-	○	条件に合致した文書単位の総件数	<code>metadataFilters</code> 省略時は <code>vectorIndex</code> 配下の全文書数
<code>documents</code>	<code>array</code>	-	○	登録済み文書のリスト	
<code>documents[].filepath</code>	<code>string</code>	-	○	登録時のファイルパス	
<code>documents[].insertDate</code>	<code>string</code>	-	○	登録時間 (UTC, ISO8601 形式)	

documents[]. url	string	-	○	ソース文書の 参照 URL	無い場合は空 文字 "" を返 す
documents[]. 任意のメタデ ータ	any	-	-	文書登録時に 任意のメタデ ータを登録し ている場合、 任意のメタデ ータを返却す る。	無い場合はキ ーごと省略

以下はレスポンスの例です。

```

1  {
2  "totalCount": 2,
3  "documents": [
4    {
5      "filepath": "file.pdf",
6      "insertDate": "2024-02-20T11:12:11Z",
7      "url": ""
8    },
9    {
10     "filepath": "test.txt",
11     "insertDate": "2024-02-20T01:12:11Z",
12     "url": "",
13     "customTag": "foo-bar"
14   }
15 ]
16 }
17

```

また、以下のcurlコマンドを用いても、文書一覧を取得できますが、メタデータによる絞込はできません。

インデックス名には、文書を登録したインデックス名を指定してください。API Keyについては管理ポータルから確認することが出来ます。

```

1  curl -X GET ¥
2  https://api.genai-api.nec-cloud.com/genai-search-api/document/listDocument?vectorIndex=<インデックス名> ¥
3  -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>"

```

取得に成功すると、以下のようにJson形式で指定したインデックスの登録済みの文章情報の一覧(ファイル名、登録時に指定したurl(指定がなければ空)、登録日時)が返却されます。

```

1  {"documents": [
2  {"filepath": "file.pdf", "url": "", "insertDate": "2024-12-09T07:16:59.147Z"},
3  {"filepath": "test.txt", "url": "https://example.com/test", "insertDate": "2024-12-10T05:20:11.304Z"}]

```

## 登録済みインデックス情報を取得する

以下のコマンド実行で作成済みのインデックス情報を取得できます。

### APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。

API Keyについては管理ポータルから確認することができます。

```
1 curl -X GET https://api.genai-api.nec-cloud.com/genai-search-api/index/listIndex ¥
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>"
```

取得に成功すると以下のようにJson形式で自分が所属しているグループが紐づいている作成済みのインデックス情報の一覧が返却されます。nameはインデックス名、document\_countは登録済み文書数、accessIdsのgroupsがインデックスに紐づいているグループのID、usersがインデックスに紐づいているユーザのID、embeddingTypeがインデックスの文書登録時に使用される埋め込みモデルの種別(現在はelasticsearchのみ)、embeddingModelが埋め込みモデル名、コメントがインデックス作成時に指定した任意の文章が表示されます。

```
1 {"indexes":[
2 {"embeddingType":"elasticsearch",{"accessIds":{"
3 "groups":["group_45b4f9ca-19d4-45d7-93eb-e469566f66dd",
4 "group_32377bc8-99d5-4137-872b-c6e97befd239"],"users":["
5 "user_fc79bb6a-b6d0-4fd3-9b45-a39525b15498",
6 "user_ad0942a4-8add-4110-882c-532b858fec1e"]}},
7 "comment":"コメント","embeddingModel":"multilingual-e5-large",
8 "name":"test-index001","document_count":1},
9 {"embeddingType":"elasticsearch",{"accessIds":{"groups":[],
10 "users":["user_aa79bb6a-b6d0-4fd3-9b45-a39525b15498",
11 "user_bb0942a4-8add-4110-882c-532b858fec1e"]}},
12 "comment":"test","embeddingModel":"multilingual-e5-large",
13 "name":"test-index002","document_count":10},
14 {"embeddingType":"elasticsearch",{"accessIds":{"groups":["
15 "group_45b4f9ca-19d4-45d7-93eb-e469566f66dd"],
16 "users":[]},"comment":"comment","embeddingModel":"multilingual-e5-large",
17 "name":"test-index003","document_count":5}]}
```

## インデックス使用量（テナント合計・インデックス別）を取得する

### インデックス使用量取得（テナント合計）

テナント内のインデックス全体の使用量（合計）と、RAG契約上の登録サイズ上限を取得します。

## 使いどころ

- ・テナント全体のインデックスの使用容量がどの程度か確認したいとき
- ・容量逼迫の監視や運用確認に使用

## APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。

API Keyについては管理ポータルから確認することができます。

```
1 curl -X GET https://api.genai-api.nec-cloud.com/genai-search-api/document/size ¥
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>"
```

取得に成功すると以下のJson形式でインデックス全体の使用量とRAG契約上の登録サイズ上限をByte単位で取得できます。

```
1 {
2   "size": {インデックス全体の使用量},
3   "maxSize": {RAG契約上の登録サイズ上限}
4 }
```

## インデックス使用量取得（インデックス別）

指定したインデックスの使用量を取得します。

## 使いどころ

- ・インデックス単位で使用容量を確認したいとき

## APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。

API Keyについては管理ポータルから確認することができます。

{index\_name}にはインデックス名を指定してください。

```
1 curl -X GET https://api.genai-api.nec-cloud.com/genai-search-api/document/size/{index_name} ¥
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>"
```

取得に成功すると以下のJson形式で特定インデックスの使用量とRAG契約上の登録サイズ上限をByte単位で取得できます。

```
1 {
2   "size": {特定インデックスの使用量},
3   "maxSize": {RAG契約上の登録サイズ上限}
4 }
```

- ① 上記のAPIが返却する使用量は、内部計算による概算値です。チャンク内容（トークン数）の偏り等により、実際の使用量と乖離する場合があります。

## 登録済みインデックス情報の変更を行う

以下のコマンド実行で登録済みのインデックス情報のうち、紐づいたグループやユーザ情報の削除とコメントの変更ができます。埋め込みモデルはインデックス作成時から変更することはできないのでご了承ください。

実行時は指定するインデックスについて、事前にAPI USERが所属しているグループ、もしくはAPI USERがインデックスに紐づいていることをご確認ください。

### APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。

API Keyについては管理ポータルから確認することが出来ます。

```
1 curl -X PUT https://api.genai-api.nec-cloud.com/genai-search-api/index/updateIndex ¥
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>" ¥
3 --data "<リクエストパラメータ>"
```

以下はリクエストパラメータの例です。vectorIndexには情報を変更したいインデックス名を指定します。

以下は任意のパラメータになります。

accessIdsのgroupsとusersは現在のインデックスに紐づけられたグループ・ユーザ情報を削除したい場合のみ以下のように指定してください。変更がない場合は指定なし("accessIds": {})で実行してください。紐づけを削除した場合、管理ポータルからグループ情報やユーザ情報を紐づけしなおす必要がありますので、注意してください。

commentはインデックスに対する説明文などを指定してください。特に変更がない場合は指定なしにしてください。

```
1 {
2   "vectorIndex": "test-index",
3   "accessIds": {"groups": [], "users": []},
4   "comment" : "テスト用インデックス"
5 }
```

変更成功すると以下のようにJson形式で指定したインデックスの名前が返却されます。

```
1 {"vectorIndex": "test-index"}
```

## 登録した文書情報を削除する

以下のコマンド実行で登録済みの文書情報をインデックスから削除できます。

実行時は指定するインデックスについて、事前にAPI USERが所属しているグループ、もしくはAPI USERがインデックスに紐づいていることをご確認ください。

## APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。インデックス名には削除したい文書が登録されているインデックス名を指定してください。

ファイル名には登録した文書のファイル名を指定して下さい。API Keyについては管理ポータルから確認することが出来ます。

```
1 curl -X DELETE ¥  
2 https://api.genai-api.nec-cloud.com/genai-search-api/document/deleteDocument?  
3 vectorIndex=<インデックス名>&filePath=<ファイル名> ¥  
4 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>"
```

文書情報の削除に成功するとnull(何もエラーなどが表示されない)になります。

## 登録した文書を一括削除する

本機能では、指定したインデックスに登録済みの文書を非同期で削除します。

また、文書登録時に付与したメタ情報（メタデータ）のKey/Value条件を指定して、検索対象を一括で指定することができます。

例えば、ファイル名や任意の追加メタデータなどのメタ情報を条件にし、条件に一致する文書を一括で検索対象とします。

## 対象となるメタデータ項目

メタデータフィルタで指定できる対象は以下です。

- source (例：ファイル名)
- url (例：参照URL)
- ユーザが任意に指定する追加メタデータ

実行時は指定するインデックスについて、事前にAPI USERが所属しているグループ、もしくはAPI USERがインデックスに紐づいていることをご確認ください。

また、メタデータ条件の指定誤り等により、削除対象が想定以上に大量にヒットする可能性があります。

そのため、誤削除を防止する目的で、削除実行前に対象件数等を確認できる 削除対象確認用のAPI(後述)を提供します。

## 削除対象確認API

本APIは、削除実行前に削除条件に一致する文書を事前に確認するためのAPIです。

本APIでは削除処理は行わず、指定した条件に一致する件数と文書情報の取得のみを行います。

### APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。API Keyは管理ポータル契約情報から確認できます。

```
1 curl -X POST https://api.genai-api.nec-cloud.com/genai-search-api/document/previewDeleteCandidates ¥  
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>" ¥  
3 --data "<リクエストパラメータ>"
```

リクエストパラメータは必須のパラメータと任意のパラメータがあります。

リクエストパラメータの各項目の詳細は以下のとおりです。

パラメータ名	型	デフォルト	必須	説明	備考
vectorIndex	string	-	△	検索対象となるベクトルストアのインデックス名を指定する。	vectorIndex、vectorIndexesのどちらかを指定する。
vectorIndexes	array[string]	-	△	検索対象となる複数のベクトルストアのインデックス名を指定する場合は配列で指定する。	vectorIndex、vectorIndexesのどちらかを指定する。
metadataAnyFilters	array[object]	-	-	登録文書のメタ情報のKey/Valueを指定して削除文書を一括で指定するための条件リスト。指定された条件のうちいずれか1つ	指定がない場合は、指定されたインデックス内のすべての文書を対象とする。

				以上に一致した文書を削除対象とする (OR 条件)。	
metadataAnyFilters[].key	string	-	○	マッチング対象とするメタデータのキー名。文書登録時のmetadataに格納されている項目名を指定する。	キーが複数階層で登録されている場合は以下のように指定する。 例：登録時 customMetadata: { "foo": { "bar": "x" } }  指定フィルタ： key="customMetadata.foo.bar", value="x"
metadataAnyFilters[].values	str, int, float, bool	-	○	比較対象となる値。文書登録時のメタデータに保存されている値を指定する。	配列にすることで、一つのkeyで複数のvalueを指定できるようにする (OR条件)
sampleSize	int	-	100	返却する文書情報の上限件数。	0~1000で指定が可能。  後述する一括削除APIでは指定不要。

**i** metadataAnyFilters は省略可能です。省略時はメタデータによる絞り込みは行わず、指定したインデックスに紐づく全文書が対象となります。

また、metadataAnyFilters を複数指定した場合、条件はORで評価されます。

以下はリクエストパラメータの例です。

```
1 {
2   "vectorIndex": "test-index",
3   "metadataAnyFilters": [
4     {
5       "key": "source",
6       "values": ["test1.txt", "test2.txt"],
7     }
8   ]
9 }
10
```

レスポンスボディの各項目の詳細は以下のとおりです。

項目	型	必須	説明	備考
matchedDocuments	int	○	削除条件 (metadataAnyFilters) に一致した文書数	
sample	array[object]	-	対象文書のサンプル (sampleSize 件、最大 1000 件)	sampleSize=0 の場合は空配列となる。

以下はレスポンスの例です。

```
1 {
2   "matchedDocuments": 2,
3   "sample": [
4     {
5       "source": "test1.txt",
6       "folder": "test-index",
7       "insertDate": "2025-09-28T04:11:47Z",
8       "url": ""
9     },
10    {
11      "source": "test2.txt",
12      "folder": "test-index",
13      "insertDate": "2025-09-29T10:05:12Z",
14      "url": "",
15      "test": 1,
16      "test2": "test2",
17      "tags": ["model", "emb1"]
18    }
19  ]
20 }
21
```

## 一括文書削除非同期API

一括削除を実行する前に、前述の削除対象確認APIを使用し、削除条件に一致する文書が意図した対象であることを確認してください。

確認後、本一括文書削除非同期APIを使用して、文書を一括で削除します。

なお、一括文書削除は非同期処理のみ提供しており、同期処理は提供していません。

## APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。API Keyは管理ポータル契約情報から確認できます。

```
1 curl -X POST https://api.genai-api.nec-cloud.com/genai-search-api/document/deleteDocumentAsync ¥
2 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>" ¥
3 --data "<リクエストパラメータ>"
```

リクエストパラメータは前述の削除対象確認APIと互換性があるため、同じリクエストパラメータを指定してください(sampleSizeのみ指定が不要です)。

- metadataAnyFiltersを指定しない場合、指定されたインデックス内のすべての文書が削除対象となります。意図しない全件削除を防ぐため、十分に注意してください。
- 削除を実行する前に、削除対象件数を事前に確認できる手段（例：削除対象確認API）を必ず実行してください。
- 削除処理を行う際は、metadataAnyFiltersを必ず指定し、削除対象を明示的に絞り込んだうえで実行してください。

レスポンスボディには以下のように削除状況を確認できるJsonのURLを返します。

```
1 {
2   "sas_url": "https://..."
3 }
4
```

## 削除状況の確認

ブラウザなどで上記のURLにアクセスすることで、以下のようなJSONファイルを参照することができます。

以下はレスポンス例です。

```
1 {
2   "id": "D0000001-delete-74da4fb8-719f-4e08-884e-d64659d86171",
3   "status": "completed",
4   "timestamp": "2025-12-16T10:30:11Z",
5   "details": {
6     "summary": {
7       "main": {
8         "total": 3,
9         "deleted": 3,
10        "failures": 0
11      }
12    },
```

```
13     "failedItems":[]
14   }
15 }
```

このJSONには削除結果に関する情報が含まれています。

## 各フィールドの説明

### • id

削除リクエストの一意のID。

### • status

削除処理の状態。以下のいずれかが設定されます。

- `received` : 受付完了
- `processing` : 削除処理中
- `completed` : 削除完了
- `partial` : 削除不足または一部失敗
- `failed` : 削除失敗

### • timestamp

ステータスが更新された日時 (UTC、ISO8601形式)。

### • details

削除結果の詳細。主に以下の情報を含みます。

#### ◦ summary

削除の集計結果。

- `total` : 対象件数
- `deleted` : 削除成功件数
- `failures` : 失敗件数

#### ◦ failedItems

削除に失敗したアイテムの情報 (配列)。失敗がない場合は空配列。

#### ◦ error ( `failed` または `partial` の場合に付与されることがあります)

エラー内容の詳細。以下の情報を含みます。

- `errorCode` : エラーコード (例: `04002305`)
- `message` : エラーメッセージ (例: `Contract information not found.`)

※ `error` は失敗時や部分失敗時にのみ含まれることがあります。成功時は出力されません。

## 作成したインデックスを削除する

以下のコマンド実行で作成済みのインデックス情報を削除できます。

実行時は指定するインデックスについて、事前にAPI USERが所属しているグループ、もしくはAPI USERがインデックスに紐づいていることをご確認ください。

## APIの実行方法

以下のcurlコマンドを用いて、APIを呼び出します。インデックス名には削除したいインデックス名を指定してください。API Keyについては管理ポータルから確認することが出来ます。

```
1 curl -X DELETE ¥  
2 https://api.genai-api.nec-cloud.com/genai-search-api/index/deleteIndex?vectorIndex=<インデックス名> ¥  
3 -H "Content-Type: application/json" -H "Authorization: Bearer <API Key>"
```

インデックスの削除に成功すると以下のように空のJsonが返ります。

```
1 {}
```