

## モデル管理APIチュートリアル はじめに

本書は、ファインチューニング済みモデルのアップロードとデプロイするチュートリアルです。

本書の対象読者は以下を想定しています。

- ・ チューニング済みモデル利用のオプションメニューを契約し、ファインチューニング済みのモデルを本サービスに登録したいSI開発者およびシステム管理者

### チュートリアルの流れ

チュートリアルは以下の流れで進みます。

1. モデルリポジトリへのアップロード
2. モデルリポジトリの一覧取得
3. モデルのデプロイ
4. デプロイ済みモデルの一覧取得
5. モデルのアンデプロイ
6. モデルリポジトリの削除

### 前提条件

ファインチューニング済みモデルをデプロイするには、オプション契約が必要です。デプロイできるモデルの数はオプション契約数に依存します。

### 事前準備

ファインチューニング済みモデルをお手元にご用意ください。

また、モデルのファイルのアップロードにAzCopyを使用します。

① LoRAモデルを利用する場合は、ベースとなるモデルを格納した場所に "lora-model" ディレクトリを作成し、LoRAモデルのファイルを "lora-model" ディレクトリに格納してください。

また、モデルの設定やパラメータを定義する config.json を、モデルを格納した場所と "lora-model" ディレクトリの両方に格納してください。

### モデルリポジトリへのアップロード

#### アップロード用URLとトークンの取得

ファインチューニング済みモデルをモデルリポジトリにアップロードします。

まず、アップロード用のURLとトークンを取得します。{モデル名}には、モデルリポジトリで使うためのモデルの名称を指定してください。

```
1 curl -X POST https://api.genai-api.nec-cloud.com/genai-modelrepos-api/modelrepos/{モデル名}/generate_upload_url ¥
2 -H "Authorization: Bearer {アクセストークン}" -H "Content-Length: 0"
```

以下のように、URLとトークンが返ってきます。

```
1 {
2   "upload_url": "https://...",
3   "token": "....."
4 }
```

トークンの有効期限は1時間です。有効期限が切れた場合は、モデルリポジトリを削除した後、アップロード用URLとトークンを取得しなおしてください。

#### ファイルのアップロード

AzCopyを使用する場合

続いてAzCopyを使用して、モデルのファイルをアップロードします。{model\_dir}にはモデルのファイルがあるディレクトリを指定してください。{upload\_url}と{token}には、上記で返ってきたURLとアクセストークンを入れてください。

```
1 azcopy copy '{model_dir}/*' '{upload_url}?{token}' --recursive
```

curlを使用する場合

curlでアップロードする場合は、各ファイルをPUTでアップロードします。{model\_dir}にはモデルのファイルがあるディレクトリを、{ファイル名}にはアップロードするファイルの名前を指定してください。{upload\_url}と{token}には、上記で返ってきたURLとアクセストークンを入れてください。また、リクエストヘッダで x-ms-blob-type: Blob 指定する必要があります。

```
1 curl -X PUT -H "x-ms-blob-type: Blob" --upload-file {model_dir}/{ファイル名} {upload_url}/{ファイル名}?{token}
```

ファイルサイズが5,000MiBを超える場合は、ファイルを4,000MiB以下のブロックに分割してアップロードする必要があります。

大きなファイル file.bin を、file.bin.001, file.bin.002, file.bin.003 に分割してアップロードする場合で説明します。

まず、以下のように分割したファイルをブロックとしてアップロードします。クエリパラメータで comp=block を指定します。

blockid はブロックを識別する一意なIDを付けて、Base64文字列にエンコードします。Base64文字列はURLエンコードする必要があります。

```
1 curl -X PUT --upload-file {model_dir}/file.bin.001 {upload_url}/file.bin?comp=block&blockid=AAAAA%3D%3D&{token}
2 curl -X PUT --upload-file {model_dir}/file.bin.002 {upload_url}/file.bin?comp=block&blockid=AAAAAB%3D%3D&{token}
3 curl -X PUT --upload-file {model_dir}/file.bin.003 {upload_url}/file.bin?comp=block&blockid=AAAAAC%3D%3D&{token}
```

次に、以下のようにブロックIDを並べたブロックリストを作成します。

```
1 <?xml version="1.0" encoding="utf-8"?>
2 <BlockList>
3   <Uncommitted>AAAAA==</Uncommitted>
4   <Uncommitted>AAAAAB==</Uncommitted>
5   <Uncommitted>AAAAAC==</Uncommitted>
6 </BlockList>
```

作成したブロックリストをPUTします。

```
1 curl -X PUT --upload-file {model_dir}/blocklist.xml {upload_url}/file.bin?comp=blocklist&{token}
```

ファイルのアップロードについての詳細は、Microsoft AzureのBLOBサービスのドキュメントを参照してください。

[Put Blob \(REST API\) - Azure Storage | Microsoft Learn](#)

[Put Block \(REST API\) - Azure Storage | Microsoft Learn](#)

[Put Block List \(REST API\) - Azure Storage | Microsoft Learn](#)

## モデルリポジトリの一覧取得

モデルリポジトリにあるモデルの一覧を取得します。アップロード用のURLを取得した時点で、モデルリポジトリに登録されます。

```
1 curl https://api.genai-api.nec-cloud.com/genai-modelreposito-api/modelreposito
2 -H "Authorization: Bearer {アクセストークン}"
```

以下のようにモデルリポジトリに登録されたモデルの情報が取得できます。

```
1 [
2   {
3     "model_name": "{モデル名}",
4     "deployed": false,
5     "aliases": []
6   },
7   ...
8 ]
```

deployed はモデルがデプロイ済みかどうかを示します。モデルがデプロイ済みの場合、デプロイ時に指定したエイリアスが aliases に入ります。また、デプロイ済みのモデルは削除できません。削除する前にアンデプロイしてください。

## モデルのデプロイ

モデルをデプロイします。以下は、curl コマンドを使用してJSONデータをAPIにPOSTする基本的な例です。

```
1 curl -X PUT https://api.genai-api.nec-cloud.com/genai-models-api/models/{エイリアス}
2 -H "Authorization: Bearer {アクセストークン}"
3 -d '{"model_name": "{モデル名}", "use_lora": true, "llm_type": "{プリセット名}"'
```

URLの {エイリアス} には、モデル名のエイリアスを指定してください。推論リクエスト時にこの名前を指定して、推論に使用するモデルを特定します。

リクエストボディの {モデル名} にはモデルリポジトリに登録したモデルの名前を指定してください。

LoRAモデルをデプロイする場合は、 `use_lora` パラメータに `true` を指定してください。 LoRAモデルを使用しない場合は、 `use_lora` パラメータに `false` を指定するか、 `use_lora` パラメータを削除します。

`llm_type` パラメータには デプロイするモデルのサイズに応じて、次のプリセット名を指定します。このパラメータを省略した場合は `g2.small` が採用されます。

プリセット名一覧：

プリセット名	モデルのサイズ
<code>g2.small</code>	cotomi-fast-v2.0 ベース
<code>g2.large</code>	cotomi-pro-v2.0 ベース
<code>g3.xlarge</code>	cotomi-v3.0 ベース

## デプロイ済みモデルの一覧取得

デプロイされたモデルを取得します。

```
1 curl https://api.genai-api.nec-cloud.com/genai-models-api/models -H "Authorization: Bearer {アクセストークン}"
```

返される一覧には、デプロイしたファインチューニング済みのモデルだけでなく、全テナント共通で提供されるNEC cotomiのモデルも含まれます。

```
1 {
2   "object": "list",
3   "data": [
4     {
5       "id": "cotomi-fast-v2.0",
6       "object": "model",
7       "created": 1234567890,
8       "owned_by": "nec",
9       "ready": true
10    },
11    {
12      "id": "cotomi-pro-v2.0",
13      "object": "model",
14      "created": 1234567890,
15      "owned_by": "nec",
16      "ready": true
17    },
18    {
19      "id": "cotomi-v3.0",
20      "object": "model",
21      "created": 1234567890,
22      "owned_by": "nec",
23      "ready": true
24    },
25    {
26      "id": "{エイリアス}",
27      "object": "model",
28      "created": 1234567890,
29      "owned_by": "{サブスクリプション名}",
30      "model_name": "{モデルリポジトリのモデル名}",
31      "ready": true,
32      "llm_type": "g2.small"
33    },
34    ...
35  ]
36 }
```

`ready` が `true` になれば、モデルにリクエストできるようになります。

## モデルのアンデプロイ

使用しないモデルをアンデプロイします。 {エイリアス} にはアンデプロイするモデルのエイリアスを指定してください。

```
1 curl -X DELETE https://api.genai-api.nec-cloud.com/genai-models-api/models/{エイリアス}
2 -H "Authorization: Bearer {アクセストークン}"
```

## モデルリポジトリの削除

モデルのファイルをモデルリポジトリから削除します。{モデル名}には、削除するモデルの名前を指定してください。デプロイ済みのモデルは削除できません。削除する前にアンデプロイしてください。

```
1 curl -X DELETE https://api.genai-api.nec-cloud.com/genai-modelrepos-api/modelrepos/{モデル名} ¥  
2 -H "Authorization: Bearer {アクセストークン}"
```