

پروژه کارشناسی:
Multimodal RAG

نام دانشجو: حسین آریان مهر
شماره دانشجویی: ۴۰۱۱۲۶۲۱۶۷
استاد: دکتر هادی صدوقی یزدی

۱ فهرست

فهرست مطالب

۲	۱ فهرست
۴	۲ نحوه <i>setup</i> پروژه و اجرای کد
۴	۱.۲ ۱ - نصب داکر
۴	۲.۲ ۲ - سورس کد پروژه
۴	۳.۲ ۳ - نصب محیط مجازی و کتابخانه‌های مورد نیاز
۴	۴.۲ ۴ - ساخت دیتابیس <i>Weaviate</i>
۴	۱.۴.۲ توضیح فایل <i>docker-compose.yml</i>
۵	۵.۲ ۵ - گرفتن API
۶	۳ ایجاد Collection در پایگاه داده <i>Weaviate</i>
۶	۱.۳ توضیح کد پایتون اتصال و پیکربندی <i>Weaviate</i>
۶	۲.۳ اتصال به <i>Weaviate</i>
۶	۳.۳ تعریف نام مجموعه (<i>Collection</i>)
۶	۴.۳ حذف مجموعه‌ی قبلی (در صورت وجود)
۶	۵.۳ ایجاد مجموعه‌ی جدید
۷	۶.۳ بستن اتصال
۷	۴ دیتاست
۸	۵ مدل <i>open_clip</i>
۸	۶ وارد کردن داده‌ها به دیتابیس
۸	۱.۶ وارد کردن داده‌های تصویری به دیتابیس
۱۱	۲.۶ وارد کردن داده‌های صوتی به دیتابیس
۱۲	۱.۲.۶ توضیح کد
۱۲	۲.۲.۶ تنظیمات مدل <i>CLIP</i>
۱۲	۳.۲.۶ مدل <i>Whisper</i>
۱۲	۴.۲.۶ بارگذاری و پردازش داده‌های صوتی
۱۳	۵.۲.۶ مرحله اول: استخراج متن از فایل‌های صوتی
۱۵	۶.۲.۶ مرحله دوم: تولید و ذخیره‌ی <i>embedding</i>
۱۹	۳.۶ فرآیند وارد کردن داده‌های متنی
۱۹	۱.۳.۶ اجرای اسکریپت
۱۹	۲.۳.۶ توضیحات جزئیات کد
۳۲	۷ بخش یک اند اصلی
۳۲	۱.۷ فایل <i>config.py</i>
۳۲	۱.۱.۷ مسیر ریشه پروژه
۳۳	۲.۱.۷ نام کالکشن پایگاه داده <i>Weaviate</i>
۳۳	۳.۱.۷ بارگذاری متغیرهای محیطی از <i>.env</i>
۳۳	۴.۱.۷ تنظیم دستگاه اجرا (<i>CPU</i> یا <i>GPU</i>)
۳۳	۵.۱.۷ پیکربندی مدل <i>CLIP</i>
۳۳	۶.۱.۷ پیکربندی مدل <i>Whisper</i>
۳۳	۷.۱.۷ پیکربندی مدل زبانی بزرگ (<i>LLM</i>)
۳۴	۲.۷ فایل <i>ai_models.py</i>
۳۴	۱.۲.۷ بارگذاری مدل <i>Whisper</i> (تبدیل گفتار به متن)
۳۴	۲.۲.۷ بارگذاری مدل <i>CLIP</i> (درک تصویر و متن)
۳۴	۳.۲.۷ بارگذاری <i>Tokenizer</i> مدل <i>CLIP</i>
۳۵	۴.۲.۷ انتقال مدل به <i>GPU</i> و حالت <i>evaluation</i>
۳۵	۳.۷ فایل <i>database.py</i>
۳۵	۱.۳.۷ توضیح دقیق عملکرد
۳۵	۴.۷ فایل <i>main.py</i>
۳۵	۱.۴.۷ کد فایل <i>main.py</i>
۳۶	۲.۴.۷ توضیح کد
۳۷	۳.۴.۷ قابلیت‌های جستجو و درخواست‌ها

۳۸	فایل <i>search_routes.py</i> و <i>utils.py</i>	۵.۷
۳۸	توضیح کلی عملکرد مسیر <i>/multimodal</i>	۱.۵.۷
۳۸	بررسی گام به گام کد	۲.۵.۷
۴۱	توابع کمکی (<i>Utility Functions</i>) در <i>utils.py</i>	۳.۵.۷
۴۶	فایل <i>llm.py</i>	۶.۷
۴۶	راه اندازی کلاینت <i>OpenAI</i>	۱.۶.۷
۴۶	تشریح تابع <i>feed_data_into_llm</i>	۲.۶.۷
۵۰	توابع کمکی	۳.۶.۷

۲ نحوه *setup* پروژه و اجرای کد

۱.۲ ۱ - نصب داکر

فایل نصبی داکر را از وبسایت رسمی آن دانلود و نصب کنید.
لینک سایت:

<https://www.docker.com/products/docker-desktop/>

پس از نصب، سیستم را ری‌استارت کنید.

۲.۲ ۲ - سورس کد پروژه

لینک پروژه در گیت:

<https://github.com/HosseinArianmehr2004/Multimodal-RAG-System>

برای دسترسی به سورس کد پروژه می‌توانید از هر یک از دو راه زیر استفاده کنید:
A - استفاده از *fork*: به آدرس داده شده در بالا بروید و در قسمت بالا سمت چپ صفحه روی فلش کنار عنوان *fork* کلیک کنید و گزینه *create a new fork* را انتخاب کرده و مراحل را پیش ببرید.
B - استفاده از *clone*: با استفاده از دستور زیر، سورس کد پروژه را در سیستم لوکال خود کلون کنید:

```
1 git clone project_link path
```

۳.۲ ۳ - نصب محیط مجازی و کتابخانه‌های موردنیاز

در پوشه‌ی پروژه (*PROJECT_ROOT*) با دستورات زیر یک محیط مجازی ایجاد کنید. یک *CMD* باز کنید و به مسیر *root* پروژه بروید. سپس دستورات زیر را اجرا کنید:

```
1 python -m venv environment_name
2 environment_name\Scripts\activate
```

پس از اکتیو شدن محیط مجازی با دستور زیر تمام کتابخانه‌های موردنیاز را نصب کنید:

```
1 pip install -r requirements.txt
```

به دلیل تعداد زیاد کتابخانه‌ها و همچنین وجود کتابخانه‌های پرحجمی مانند *pytorch* در بین آن‌ها، زمان نصب ممکن است کمی طول بکشد.

۴.۲ ۴ - ساخت دیتابیس *Weaviate*

در ابتدا نرم‌افزار *Docker* را در ویندوز اجرا کنید. سپس یک *CMD* باز کرده و به مسیر *root* پروژه بروید. برای فعال‌سازی سرویس *Weaviate* در داکر، دستور زیر را اجرا کنید:

```
1 docker-compose up -d
```

این دستور فایل *docker-compose.yml* را اجرا کرده و سرویس *Weaviate* را روی داکر راه‌اندازی می‌کند.

۱.۴.۲ توضیح فایل *docker-compose.yml*

فایل *docker-compose.yml* وظیفه دارد نحوه‌ی استقرار و اجرای سرویس‌های مختلف پروژه را در محیط *Docker* مشخص کند. در این پروژه، تنها یک سرویس با نام *Weaviate* تعریف شده است. محتوای این فایل به شرح زیر است:

```
1 version: '3.4'
2 services:
3   weaviate:
4     image: semitechnologies/weaviate:latest
5     restart: always
6     ports:
7       - "8080:8080"
8       - "50051:50051"
9     environment:
10       QUERY_DEFAULTS_LIMIT: 25
11       AUTHENTICATION_ANONYMOUS_ACCESS_ENABLED: 'true'
12       PERSISTENCE_DATA_PATH: "./data"
13       DISABLE_MODULES: 'all'
```

توضیحات بخش‌های مختلف:

- **version**: نسخه‌ی نگارش فایل *docker-compose* را مشخص می‌کند. در نسخه‌های جدید *Docker* این گزینه اختیاری است و می‌توان آن را حذف کرد.
- **services**: لیست سرویس‌هایی است که باید توسط *Docker* اجرا شوند. در اینجا تنها سرویس *weaviate* تعریف شده است.
- **weaviate**: نام سرویس مربوط به پایگاه‌داده‌ی برداری (*Weaviate*) است.
- **image**: مشخص می‌کند که سرویس از تصویر (*Image*) رسمی *Weaviate* با آخرین نسخه (*latest*) استفاده می‌کند. این تصویر از مخزن *Docker Hub* دریافت می‌شود.
- **restart**: با مقدار *always* تعیین می‌شود که در صورت توقف یا بروز خطا، کانتینر به‌صورت خودکار مجدداً راه‌اندازی شود.
- **ports**: نگاشت پورت‌های داخلی کانتینر به پورت‌های میزبان را تعیین می‌کند:
 - پورت 8080 برای دسترسی به رابط *HTTP* و *API* اصلی *Weaviate*.
 - پورت 50051 برای ارتباطات *gRPC*.
- **environment**: شامل متغیرهای محیطی است که پارامترهای اجرایی *Weaviate* را پیکربندی می‌کنند:
 - *QUERY_DEFAULTS_LIMIT: 25* — تعداد نتایج پیش‌فرض در هر Query را ۲۵ تعیین می‌کند.
 - *AUTHENTICATION_ANONYMOUS_ACCESS_ENABLED: 'true'* — امکان دسترسی ناشناس به *Weaviate* را فعال می‌سازد (بدون نیاز به احراز هویت).
 - *PERSISTENCE_DATA_PATH: "/data"* — مسیر ذخیره‌سازی داده‌های پایدار را مشخص می‌کند تا داده‌ها حتی پس از توقف کانتینر حذف نشوند.
 - *DISABLE_MODULES: 'all'* — تمام ماژول‌های اختیاری *Weaviate* (از جمله مدل‌های بردارساز) را غیرفعال می‌کند و تنها هسته‌ی اصلی پایگاه‌داده اجرا می‌شود.

نتیجه‌ی اجرا: با اجرای دستور زیر:

```
docker-compose up -d
```

سرویس تحت عنوان *Weaviate* در محیط *Docker* اجرا می‌شود که از طریق آدرس‌های زیر در دسترس است:

● *HTTP API: http://localhost:8080*

● *gRPC API: localhost:50051*

این تنظیمات باعث می‌شود *Weaviate* به‌صورت پایدار و مستقل در بستر *Docker* اجرا گردد، داده‌ها در مسیر محلی ذخیره شوند و امکان ارتباط از طریق پورت‌های مشخص‌شده فراهم گردد.

۵.۲ - گرفتن API

برای استفاده از مدل‌های *LLM* نیاز به *API* آن‌ها دارید. در این پروژه، *API*‌های مدنظر از *OpenRouter* گرفته شده‌اند. برای گرفتن *API* به طریق زیر عمل کنید:

۱. به وب‌سایت رسمی *OpenRouter* بروید: <https://openrouter.ai/>

۲. در صورتی که حساب کاربری ندارید، ابتدا یک حساب کاربری بسازید و لاگین کنید.

۳. روی علامت منو در بالا سمت راست صفحه کلیک کرده و گزینه‌ی *keys* را انتخاب کنید.

۴. در صفحه‌ی مورد نظر، روی دکمه‌ی *Create API Key* کلیک کنید.

۵. یک نام برای کلید خود انتخاب کرده و گزینه‌ی *Create* را بزنید.

۶. سپس یک *pop-up* به شما نمایش داده می‌شود. از همین قسمت کلید خود را کپی و در یک فایل ذخیره کنید.

توجه کنید که پس از بستن *pop-up*، نمی‌توانید به محتوای کلید دست یابید. سپس باید کلید خود را در متغیرهای محیطی پروژه ذخیره کنید تا در کد از آن استفاده شود.

برای این کار، به پوشه‌ی پروژه رفته و در آنجا فایلی با نام *env*. بسازید و خط زیر را در آن کپی کنید. سپس کلید خود را در قسمت *your_key* قرار دهید:

```
LLM_API_KEY=your_key
```

در پایان، فایل را ذخیره کنید.

۳ ایجاد Collection در پایگاه داده Weaviate

برای ذخیره سازی داده ها در پایگاه داده باید یک *collection* ایجاد کنید. برای این کار ابتدا یک *cmd* باز کرده و به مسیر *root* پروژه بروید، سپس اسکریپت *create_database_collection.py* را به صورت زیر اجرا نمایید:

```
python import_data_into_database/create_database_collection.py
```

۱.۳ توضیح کد پایتون اتصال و پیکربندی Weaviate

کد زیر به منظور اتصال به پایگاه داده ی برداری *Weaviate*، ایجاد یک مجموعه (*Collection*) و تعریف ساختار داده ها در آن نوشته شده است. این کد از کتابخانه ی رسمی *weaviate-client* در پایتون استفاده می کند.

```
1 import weaviate
2 import weaviate.classes as wvc
```

دو ماژول اصلی کتابخانه *Weaviate* فراخوانی می شوند:

- *weaviate*: برای مدیریت اتصال و انجام عملیات اصلی.
- *weaviate.classes*: برای تعریف تنظیمات و ساختار مجموعه ها (*Collections*) و ویژگی ها (*Properties*).

۲.۳ اتصال به Weaviate

```
1 try:
2     client = weaviate.connect_to_local()
3     print("Successfully connected to Weaviate!")
4 except Exception as e:
5     print(f"Failed to connect to Weaviate: {e}")
6     exit()
```

در این بخش، با استفاده از تابع *connect_to_local()* یک اتصال به سرویس محلی *Weaviate* برقرار می شود (یعنی همان سرویسی که با *Docker* روی *localhost:8080* اجرا شده است).

۳.۳ تعریف نام مجموعه (Collection)

```
1 collection_name = "Multimodal_Collection"
```

در این خط، نام مجموعه ای که قرار است در پایگاه داده ایجاد شود مشخص می گردد. در *Weaviate*، مجموعه ها مشابه «جدول» در پایگاه داده های رابطه ای هستند.

۴.۳ حذف مجموعه ی قبلی (در صورت وجود)

```
1 if client.collections.exists(collection_name):
2     print(f"Collection '{collection_name}' already exists. Deleting it.")
3     client.collections.delete(collection_name)
```

در این بخش ابتدا بررسی می شود که آیا مجموعه ای با همین نام از قبل وجود دارد یا خیر. اگر وجود داشته باشد، برای جلوگیری از تکرار و تعارض، ابتدا حذف می شود.

۵.۳ ایجاد مجموعه ی جدید

```
1 print(f"Creating collection '{collection_name}'...")
2 my_collection = client.collections.create(
3     name=collection_name,
4     properties=[
5         wvc.config.Property(name="modality", data_type=wvc.config.DataType.TEXT),
6         wvc.config.Property(name="content", data_type=wvc.config.DataType.TEXT),
7         wvc.config.Property(name="contentId", data_type=wvc.config.DataType.TEXT),
8         wvc.config.Property(name="filePath", data_type=wvc.config.DataType.TEXT),
9     ]
10 )
11 print(f"Collection '{collection_name}' created successfully.")
```

در این مرحله، مجموعه‌ای جدید با نام *Multimodal_Collection* ایجاد می‌شود. پارامتر *properties* مشخص‌کننده‌ی فیلدهای درون این مجموعه است. هر ویژگی با استفاده از کلاس *wvc.config.Property* تعریف می‌شود و نوع داده‌ی آن نیز با *wvc.config.DataType* تعیین می‌گردد.

نام ویژگی	نوع داده	توضیح
modality	TEXT	نوع داده‌ی ورودی (مثلاً متن، تصویر یا صوت)
content	TEXT	محتوای اصلی داده
contentId	TEXT	شناسه‌ی منحصر به فرد هر داده یا سند
filePath	TEXT	مسیر فایل اصلی در سیستم محلی یا فضای ذخیره‌سازی

به این ترتیب، ساختار مجموعه به گونه‌ای طراحی شده که از داده‌های چندرسانه‌ای (متن، تصویر و صوت) پشتیبانی کند. در این مجموعه، هر داده شامل فیلدهای بالا است:

داده‌های متنی

```
modality = "text"
content = Original text
contentId = a unique id
filePath = ""
```

داده‌های تصویری

```
modality = "image"
content = ""
contentId = a unique id
filePath = "/content/image_dataset/{file_name}"
```

داده‌های صوتی

```
modality = "audio"
content = Audio transcription
contentId = a unique id
filePath = "/content/audio_dataset/{file_name}"
```

۶.۳ بستن اتصال

```
client.close()
```

پس از اتمام عملیات، اتصال به پایگاه داده بسته می‌شود تا منابع آزاد شوند.

۴ دیتاست

دیتاست‌ها در فضای *Google Drive* ذخیره شده‌اند و از طریق لینک‌های زیر قابل دانلود هستند:

- دیتاست تصاویر:

<https://drive.google.com/file/d/1a68xCnTnnRPrBSM7bHSRzHwLBp0oLWuL/view?usp=sharing>

- دیتاست فایل‌های صوتی:

https://drive.google.com/file/d/1BS72l_5Z6wvGiq9Sn8zuV4-IfwxLRvR5/view?usp=sharing

- دیتاست متن:

<https://drive.google.com/file/d/1RMjyguno7iLfUtlW9gumdwrgbB9ht252/view?usp=sharing>

در پوشه‌ی پروژه، یک پوشه‌ی جدید به نام *content* ایجاد کنید. در این پوشه، سه زیرپوشه با نام‌های زیر بسازید و داده‌های مربوطه را در هر کدام قرار دهید:

- *audio_dataset*: شامل تمام فایل‌های صوتی

- *image_dataset*: شامل تمام تصاویر

● *text_dataset*: شامل فایل *CSV* داده‌های متنی

در صورتی که قصد دارید از دیتاست‌ها یا داده‌های خود استفاده کنید، روال به همان شکل بالا است. کافی است تمام تصاویر خود را در پوشه‌ی *image_dataset* و تمام فایل‌های صوتی را در پوشه‌ی *audio_dataset* قرار دهید. برای داده‌های متنی، باید یک فایل *CSV* با ساختار زیر ایجاد نمایید: متن‌های طولانی، فایل‌های *PDF* و سایر متون را به رشته‌های (*string*) جداگانه تقسیم کنید. سپس فایلی با فرمت *CSV* ایجاد کرده و در آن ستونی به نام *text* قرار دهید. تمام داده‌های متنی (رشته‌ها) را در این ستون وارد کنید.

```
1 text
2 data 1
3 data 2
4 data 3
5 ...
```

داده‌های متنی و صوتی شما می‌توانند به زبان‌های فارسی یا انگلیسی باشند.

۵ مدل *open_clip*

برای استفاده از مدل *open_clip* باید وزن‌های آن را دانلود کنید. وزن‌ها را می‌توانید از لینک زیر در گوگل درایو دریافت کنید:

<https://drive.google.com/file/d/1czuHrdWylmzYHA1Rh5NQJsnFDB9BSp7Y/view?usp=sharing>

پس از دانلود، فایل را در مسیر زیر قرار دهید:

`PROJECT_ROOT/open_clip_weights/ViT-B-32-openai/`

۶ وارد کردن داده‌ها به دیتابیس

۱.۶ وارد کردن داده‌های تصویری به دیتابیس

پس از قرار دادن تمام تصاویر در پوشه مربوطه، یک *CMD* باز کنید و به مسیر پروژه بروید. سپس با دستور زیر اسکریپت پایتون با نام *images.py* را اجرا کنید:

```
1 python import_data_into_database/import_local_data/images.py
```

این اسکریپت تمام داده‌های داخل پوشه *PROJECT_ROOT/content/image_dataset* را بارگذاری میکند، برای هر کدام یک *embedding* تولید میکند و سپس آن *embedding* را به همراه مسیر فایل و چند اطلاعات دیگر در دیتابیس ذخیره میکند

توضیح کد

این اسکریپت در ابتدا با استفاده از کد زیر به دیتابیس و *collection* ای که از قبل ساخته‌ایم متصل می‌شود:

```
1 try:
2     weaviate_client = weaviate.connect_to_local()
3     print(" Connected to Weaviate")
4     collection = weaviate_client.collections.get("Multimodal_Collection")
5 except Exception as e:
6     print(f" Weaviate connection failed: {e}")
7     weaviate_client = None
```

در این بخش، با استفاده از تابع *connect_to_local()* یک اتصال به سرویس محلی *Weaviate* برقرار می‌شود (همان سرویسی که با *Docker* روی *localhost:8080* اجرا شده است).

تنظیمات مدل *CLIP*

کد زیر تنظیمات مدل *CLIP* را انجام می‌دهد؛ این مدل برای تولید *embedding* داده‌ها به کار می‌رود:

```
1 MODEL_NAME = "ViT-B-32"
2 PRETRAINED_LOCAL_PATH = (
3     "./open_clip_weights/ViT-B-32-openai/open_clip_model.safetensors"
4 )
5 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
6
7 clip_model, _, preprocess = open_clip.create_model_and_transforms(
8     MODEL_NAME, pretrained=PRETRAINED_LOCAL_PATH
9 )
```



```

10 tokenizer = open_clip.get_tokenizer(MODEL_NAME)
11 clip_model.to(DEVICE).eval()

```

توضیح پارامترها:

- **MODEL_NAME** = "ViT-B-32": نام معماری مدل CLIP (Vision Transformer Base with patch size = 32×32).
- **PRETRAINED_LOCAL_PATH**: مسیر فایل وزن‌های از پیش‌آموزش داده‌شده‌ی مدل.
- **DEVICE**: بررسی وجود GPU (CUDA)، در صورت نبود GPU از CPU استفاده می‌شود.
- **preprocess_clip_model**: بارگذاری مدل و تابع پیش‌پردازش تصاویر.
- **tokenizer**: توکنایزر مخصوص همان مدل.
- **clip_model.to(DEVICE).eval()**: انتقال مدل به GPU یا CPU و قرار دادن در حالت ارزیابی.

بارگذاری و پردازش تصاویر

```

1 if __name__ == "__main__":
2     image_folder = "./content/image_dataset"
3     process_images(image_folder)
4     weaviate_client.close()

```

توضیح کد:

- مسیر پوشه حاوی تصاویر مشخص شده است.
- تابع **process_images** برای تمام تصاویر داخل پوشه مربوطه **embedding** تولید می‌کند و آنها را در دیتابیس ذخیره می‌کند.

تابع **process_images**

```

1 def process_images(image_folder: str, max_workers: int = 8):
2     """Process and store image embeddings in parallel."""
3     images = [
4         file_path
5         for file_path in os.listdir(image_folder)
6         if file_path.lower().endswith((".jpg", ".jpeg", ".png"))
7     ]
8     with ThreadPoolExecutor(max_workers=max_workers) as executor:
9         futures = [
10             executor.submit(
11                 store_image_item, os.path.splitext(os.path.basename(img_name))[0],
12                 img_name
13             )
14             for img_name in images
15         ]
16     for _ in tqdm(as_completed(futures), total=len(futures), desc="Processing images"):
17         pass

```

توضیح تابع

• ورودی‌ها:

- **image_folder**: مسیر پوشه‌ای که تصاویر داخل آن هستند.
- **max_workers**: حداکثر تعداد **thread**هایی که همزمان می‌توانند کار کنند (پیش‌فرض ۸).

• جمع‌آوری مسیر تصاویر:

- با **os.listdir(image_folder)** تمام فایل‌های داخل پوشه گرفته می‌شوند.
- با شرط **endswith** فقط فایل‌هایی که پسوند **.jpg**، **.jpeg** یا **.png** دارند انتخاب می‌شوند.

– نتیجه یک لیست از نام فایل‌ها است که به images ذخیره می‌شود.

- پردازش موازی با ThreadPoolExecutor:

– ThreadPoolExecutor یک thread pool ایجاد می‌کند تا چند تصویر همزمان پردازش شوند.

– executor.submit(func, arg1, arg2, ...) یک کار (task) برای اجرای تابع func با آرگومان‌های مشخص ایجاد می‌کند.

– برای هر تصویر:

* نام فایل بدون پسوند (os.path.splitext(os.path.basename(img_name))[0]) به عنوان item_id استفاده می‌شود.

* مسیر تصویر به store_image_item داده می‌شود تا embedding آن ذخیره شود.

– نتیجه یک لیست از futures است. هر future نماینده نتیجه یک کار در حال اجرا است.

- نمایش پیشرفت پردازش:

– as_completed(futures) یک generator برمی‌گرداند که هر بار یکی از کارها تمام شود، آن را yield می‌کند.

– tqdm(..., desc="Processing images") یک نوار پیشرفت (progress bar) نمایش می‌دهد.

تابع store_image_item

در تابع process_images از تابع زیر استفاده شده است

```
1 def store_image_item(item_id: str, img_name: str):
2     """Store image embedding and metadata in Weaviate."""
3     embedding = get_embedding("image", img_name)
4     relative_path = f"/content/image_dataset/{img_name}"
5     properties = {
6         "contentId": item_id,
7         "modality": "image",
8         "filePath": relative_path,
9         "content": "",
10    }
11    collection.data.insert(properties=properties, vector=embedding.tolist())
```

توضیح پارامترها و مراحل:

- پارامترها:

– item_id: شناسه یکتا برای هر تصویر (معمولاً نام فایل بدون پسوند).

– img_name: نام فایل تصویر یا مسیر نسبی آن.

- تولید embedding تصویر:

– با استفاده از تابع get_embedding("image", img_name)، تصویر به یک بردار عددی (vector) تبدیل می‌شود.

– این embedding برای جستجوی شباهت و یادگیری ماشین استفاده می‌شود.

- تعیین مسیر نسبی تصویر:

– relative_path = f"/content/image_dataset/{img_name}" برای ذخیره در metadata ساخته می‌شود.

- آماده‌سازی metadata:

– properties شامل ویژگی‌های تصویر است:

* contentId: شناسه یکتا تصویر

* modality: نوع داده، اینجا "image"

* filePath: مسیر تصویر

* content: خالی برای تصاویر

- ذخیره در Weaviate:

– collection.data.insert(properties=properties, vector=embedding.tolist()) داده‌ها

و embedding تصویر را در دیتابیس ذخیره می‌کند.

– این کار تصویر و بردار ویژگی آن را برای جستجوی مشابهت و بازیابی سریع آماده می‌کند.

تابع get_embedding

برای تولید embedding برای تصاویر از تابع زیر استفاده شده است

```
1 def get_embedding(modality: str, image_name: Union[str, None]) -> np.ndarray:
2     mod = modality.lower()
3     if mod == "image":
4         if not isinstance(image_name, str):
5             raise ValueError("`image_name` must be a string.")
6
7         img_path = f"./content/image_dataset/{image_name}"
8         img = Image.open(img_path).convert("RGB")
9         x = preprocess(img).unsqueeze(0).to(DEVICE)
10
11         with torch.no_grad():
12             features = clip_model.encode_image(x)
13             features = features / features.norm(dim=-1, keepdim=True)
14         return features.detach().cpu().numpy().reshape(-1)
15     else:
16         raise ValueError("`modality` must be 'image'.")
```

توضیح تابع:

• ورودی‌ها:

— modality: نوع داده

— image_name: نام فایل تصویر (یا None اگر ورودی نامعتبر باشد).

• بارگذاری و آماده‌سازی تصویر:

— مسیر تصویر ساخته می‌شود: `img_path = f"./content/image_dataset/{image_name}"`

— تصویر با PIL باز و به RGB تبدیل می‌شود.

— با تابع `preprocess` آماده‌سازی می‌شود (تغییر اندازه، نرمال‌سازی و تبدیل به *Tensor*).

— `unsqueeze(0)` برای اضافه کردن بعد batch و `to(DEVICE)` برای انتقال به CPU/GPU.

• تولید embedding با CLIP:

— با `torch.no_grad()` محاسبه گرادیان غیرفعال می‌شود (صرفه‌جویی در حافظه و سرعت).

— `clip_model.encode_image(x)` تصویر را به یک بردار ویژگی (embedding) تبدیل می‌کند.

— نرمال‌سازی: `features / features.norm(dim=-1, keepdim=True)` تا طول بردار برابر ۱ شود (مناسب جستجوی شباهت cosine).

• تبدیل به NumPy و بازگرداندن:

— `detach()` از گراف محاسباتی جدا می‌کند.

— `cpu()` بردار را روی CPU قرار می‌دهد.

— `numpy()` به آرایه NumPy تبدیل می‌کند.

— `reshape(-1)` بردار را یک‌بعدی می‌کند و باز می‌گرداند.

۲.۶ وارد کردن داده‌های صوتی به دیتابیس

در این بخش، فرآیند وارد کردن داده‌های صوتی به پایگاه داده توضیح داده می‌شود. در ابتدا، پس از قرار دادن تمامی فایل‌های صوتی در پوشه‌ی مربوطه، یک پنجره‌ی `cmd` باز کرده و به مسیر پروژه بروید. سپس با استفاده از دستور زیر، اسکریپت پایتون با نام `audio.py` اجرا می‌شود:

```
python import_data_into_database/import_local_data/audio.py
```

این اسکریپت تمام فایل‌های صوتی داخل پوشه `PROJECT_ROOT/content/audio_dataset` را خوانده و با استفاده از مدل Whisper متن آنها را به دست می‌آورد و همچنین متن‌های غیرانگلیسی را به فارسی ترجمه می‌کند. سپس برای متن‌های به دست آمده با استفاده از مدل CLIP یک embedding تولید می‌کند و در نهایت این embedding را به همراه مسیر فایل صوتی و همچنین متن به دست آمده در دیتابیس ذخیره می‌کند.

۱.۲.۶ توضیح کد

در ابتدای اجرای این اسکریپت، اتصال به پایگاه داده‌ی *Weaviate* و کالکشن (*collection*) از پیش ساخته شده برقرار می‌شود. این کار توسط کد زیر انجام می‌گیرد:

```
1 try:
2     weaviate_client = weaviate.connect_to_local()
3     print(" Connected to Weaviate")
4     collection = weaviate_client.collections.get("Multimodal_Collection")
5 except Exception as e:
6     print(f" Weaviate connection failed: {e}")
7     weaviate_client = None
```

در این بخش، تابع *connect_to_local()* اتصال به سرویس محلی *Weaviate* را برقرار می‌کند (همان سرویسی که با *Docker* روی *localhost:8080* اجرا شده است).

۲.۲.۶ تنظیمات مدل CLIP

در ادامه، تنظیمات مربوط به مدل *CLIP* انجام می‌شود. این مدل برای تولید *embedding* داده‌ها به کار می‌رود:

```
1 MODEL_NAME = "ViT-B-32"
2 PRETRAINED_LOCAL_PATH = (
3     "./open_clip_weights/ViT-B-32-openai/open_clip_model.safetensors"
4 )
5 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
6
7 clip_model, _, preprocess = open_clip.create_model_and_transforms(
8     MODEL_NAME, pretrained=PRETRAINED_LOCAL_PATH
9 )
10 tokenizer = open_clip.get_tokenizer(MODEL_NAME)
11 clip_model.to(DEVICE).eval()
```

- *MODEL_NAME*: نام معماری مدل *CLIP* را مشخص می‌کند (مدل *Vision Transformer Base* با اندازه‌ی پیچ 32×32).
- *PRETRAINED_LOCAL_PATH*: مسیر فایل وزن‌های از پیش آموزش داده شده‌ی مدل را تعیین می‌کند تا از آن بارگذاری شود.
- *DEVICE*: بررسی می‌کند که آیا *GPU (CUDA)* در دسترس است یا خیر؛ در صورت وجود از *GPU* و در غیر این صورت از *CPU* استفاده می‌شود.
- *create_model_and_transforms*: مدل *CLIP* و توابع پیش پردازش آن را بارگذاری می‌کند.
- *tokenizer*: توکنایزر مخصوص مدل را دریافت می‌کند تا متون به توکن‌های عددی تبدیل شوند.
- *clip_model.to(DEVICE).eval()*: مدل را به دستگاه مربوطه منتقل کرده و در حالت ارزیابی قرار می‌دهد.

۳.۲.۶ مدل Whisper

در این بخش از مدل *Whisper* برای استخراج متن از فایل‌های صوتی استفاده می‌شود:

```
1 MODEL_TYPE = "small"
2 whisper_model = whisper.load_model(MODEL_TYPE, device=DEVICE)
```

مدل *Whisper* به کمک پارامتر *MODEL_TYPE* مقداردهی می‌شود. این پارامتر می‌تواند یکی از مقادیر *small*، *medium* یا *large* باشد. مدل‌های *medium* و *large* تنها بر روی *GPU* قابل اجرا هستند، در حالی که مدل *small* بر روی *CPU* نیز قابل بارگذاری است، هرچند دقت پایین‌تری دارد. همچنین، مدل *small* برای فایل‌های صوتی زبان انگلیسی عملکرد مناسبی دارد ولی برای زبان فارسی ضعیف‌تر است؛ بنابراین در صورت نیاز به پشتیبانی از زبان فارسی، استفاده از مدل‌های *medium* یا *large* توصیه می‌شود.

۴.۲.۶ بارگذاری و پردازش داده‌های صوتی

در ادامه، داده‌های صوتی بارگذاری و پردازش می‌شوند:

```

1 if __name__ == "__main__":
2     audio_folder = "./content/audio_dataset"
3     audio_json = f"{audio_folder}/audio_metadata.json"
4
5     transcribe_audios(audio_folder, audio_json)
6     process_audio_json(audio_json)
7     weaviate_client.close()

```

مسیر فایل‌ها به صورت خودکار تنظیم می‌شود و نیازی به تغییر ندارد. پردازش داده‌های صوتی در دو مرحله انجام می‌شود:

۱. استخراج متن (*Transcription*) از فایل‌های صوتی؛

۲. تولید *embedding* از متون و ذخیره در پایگاه داده.

۵.۲.۶ مرحله اول: استخراج متن از فایل‌های صوتی

این کار با تابع زیر انجام می‌شود:

```

1 def transcribe_audios(folder, out_json="results.json"):
2     if not os.path.isdir(folder):
3         raise ValueError(f"No folder: {folder}")
4
5     files = sorted(
6         [
7             os.path.join(folder, file_path)
8             for file_path in os.listdir(folder)
9             if file_path.lower().endswith((".wav", ".mp3", ".m4a"))
10        ]
11    )
12
13    results = [
14        {
15            "id": i + 1,
16            "filename": os.path.basename(audio_path),
17            "audio_path": audio_path,
18            "transcription": transcribe_to_english(audio_path),
19        }
20        for i, audio_path in enumerate(tqdm(files, desc="Processing"))
21    ]
22
23    with open(out_json, "w", encoding="utf-8") as f:
24        json.dump(results, f, ensure_ascii=False, indent=2)
25    print(f" Saved {len(results)} items to {out_json}")
26    return results

```

تعریف تابع و ورودی‌ها

```

1 def transcribe_audios(folder, out_json="results.json"):

```

- *folder*: مسیر پوشه‌ای که فایل‌های صوتی در آن قرار دارند.
- *out_json*: نام فایل خروجی با فرمت *JSON* (به صورت پیش فرض برابر با "results.json").

جمع‌آوری فایل‌های صوتی

در گام نخست، تابع تمامی فایل‌های صوتی موجود در پوشه‌ی مشخص شده را شناسایی و آماده‌ی پردازش می‌کند.

```

1 files = sorted(
2     [
3         os.path.join(folder, file_path)
4         for file_path in os.listdir(folder)
5         if file_path.lower().endswith((".wav", ".mp3", ".m4a"))
6     ]
7 )

```

در این بخش:

- تمامی فایل‌های موجود در پوشه بررسی می‌شوند.
- تنها فایل‌هایی که پسوند آن‌ها *wav*، *mp3* یا *m4a* باشد انتخاب می‌گردند.

- مسیر کامل هر فایل با استفاده از تابع `os.path.join` ساخته می‌شود.

- سپس فهرست فایل‌ها با دستور `sorted` مرتب می‌شود.

در نهایت، متغیر `files` شامل فهرستی از مسیر کامل تمام فایل‌های صوتی موجود در پوشه است.

ایجاد خروجی برای هر فایل

در این مرحله، برای هر فایل صوتی موجود در فهرست، خروجی متنی (متن استخراج‌شده از گفتار) تولید و در قالب یک ساختار

داده‌ی دیکشنری ذخیره می‌گردد.

```
1 results = [  
2     {  
3         "id": i + 1,  
4         "filename": os.path.basename(audio_path),  
5         "audio_path": audio_path,  
6         "transcription": transcribe_to_english(audio_path),  
7     }  
8     for i, audio_path in enumerate(tqdm(files, desc="Processing"))  
9 ]
```

در این بخش:

- از کتابخانه‌ی `tqdm` برای نمایش نوار پیشرفت استفاده می‌شود.

- برای هر فایل صوتی:

- `id`: شماره‌ی ترتیبی فایل (از عدد ۱ آغاز می‌شود).

- `filename`: نام فایل بدون مسیر.

- `audio_path`: مسیر کامل فایل.

- `transcription`: نتیجه‌ی حاصل از اجرای تابع `transcribe_to_english(audio_path)` که متن انگلیسی استخراج‌شده

از گفتار در فایل است.

در پایان، خروجی این بخش شامل فهرستی از دیکشنری‌هاست که هر کدام نمایانگر اطلاعات و متن استخراج‌شده از یک فایل صوتی می‌باشد.

ذخیره در فایل JSON

در این مرحله، داده‌های استخراج‌شده در قالب فایل `JSON` ذخیره می‌شوند.

```
1 with open(out_json, "w", encoding="utf-8") as f:  
2     json.dump(results, f, ensure_ascii=False, indent=2)
```

توضیحات:

- فایل خروجی `JSON` باز می‌شود و محتوای نتایج در آن نوشته می‌شود.

- پارامتر `ensure_ascii=False` تضمین می‌کند که کاراکترهای غیرانگلیسی (مانند فارسی) به‌درستی ذخیره شوند.

- مقدار `indent=2` برای زیبایی و خوانایی ساختار فایل خروجی استفاده می‌گردد.

تابع `transcribe_to_english`

برای استخراج متن از هر فایل صوتی، از تابع `transcribe_to_english` استفاده می‌شود. این تابع وظیفه دارد یک فایل صوتی

را دریافت کرده و متن آن را به زبان انگلیسی برگرداند، صرف‌نظر از اینکه زبان اصلی گفتار چه باشد.

```
1 def transcribe_to_english(path):  
2     res = whisper_model.transcribe(path, task="transcribe")  
3     lang, text = res.get("language", ""), res.get("text", "").strip()  
4     if lang.startswith("en"):  
5         return text  
6     else:  
7         return (  
8             whisper_model.transcribe(path, task="translate", language="en")  
9             .get("text", "")  
10            .strip()  
11        )
```

ورودی

- `path`: مسیر فایل صوتی ورودی.

استخراج گفتار (Transcription)

```
1 res = whisper_model.transcribe(path, task="transcribe")
```

در این مرحله از مدل *Whisper* (مدل تشخیص گفتار از *OpenAI*) استفاده می‌شود. پارامتر `task="transcribe"` به مدل اعلام می‌کند که تنها گفتار را به همان زبان اصلی به متن تبدیل کند، بدون ترجمه. خروجی حاصل در متغیر `res` ذخیره می‌شود که ساختاری مشابه زیر دارد:

```
1 {
2   "text": "Bonjour tout le monde",
3   "language": "fr"
4 }
```

استخراج زبان و متن

```
1 lang, text = res.get("language", ""), res.get("text", "").strip()
```

در این مرحله:

- `lang`: زبان تشخیص داده شده توسط مدل (مانند "en"، "fr"، "fa"، "de" و ...).
- `text`: متن استخراج شده از گفتار (پس از حذف فاصله‌های اضافی).

بررسی زبان انگلیسی

```
1 if lang.startswith("en"):
2     return text
```

اگر زبان تشخیص داده شده انگلیسی باشد (یعنی مقدار `lang` با "en" شروع شود)، همان متن اولیه مستقیماً بازگردانده می‌شود. ترجمه در صورت غیرانگلیسی بودن زبان

```
1 else:
2     return (
3         whisper_model.transcribe(path, task="translate", language="en")
4         .get("text", "")
5         .strip()
6     )
```

در صورتی که زبان اصلی غیرانگلیسی باشد:

- مدل *Whisper* مجدداً اجرا می‌شود، اما این بار با پارامترهای:
 - `task="translate"`: به این معنا که مدل همزمان با تشخیص گفتار، ترجمه به انگلیسی را نیز انجام دهد.
 - `language="en"`: خروجی نهایی حتماً به زبان انگلیسی باشد.
- در انتها، تنها مقدار مربوط به کلید "text" استخراج شده و فاصله‌های اضافی آن حذف می‌شود.

خروجی نهایی

تابع همواره یک رشته متنی (*string*) بازمی‌گرداند که نسخه‌ی انگلیسی گفتار موجود در فایل صوتی است، صرف‌نظر از اینکه زبان اصلی صوت چه بوده است. در پایان این مرحله، یک فایل *JSON* شامل فراداده‌ها (*metadata*) و متن‌های استخراج شده در همان پوشه‌ای که فایل‌های صوتی قرار دارند تولید می‌شود. به این ترتیب، برای هر داده‌ی صوتی، یک متن معادل انگلیسی نیز در اختیار خواهیم داشت.

۶.۲.۶ مرحله دوم: تولید و ذخیره‌ی *embedding*

در این مرحله، با استفاده از تابع `process_audio_json`، فایل *JSON* مرحله‌ی قبل خوانده می‌شود و برای هر داده صوتی *embedding* تولید می‌گردد:

```
1 def process_audio_json(json_path: str, max_workers: int = 8):
2     with open(json_path, "r") as f:
3         metadata = json.load(f)
4
5     with ThreadPoolExecutor(max_workers=max_workers) as executor:
6         futures = []
7         for item in metadata:
8             file_path = item.get("audio_path", "").strip()
9             transcription = item.get("transcription", "").strip()
10            if not (file_path and transcription):
```

```

11         continue
12         content_id = str(item.get("id", os.path.basename(file_path)))
13         futures.append(
14             executor.submit(
15                 store_audio_item,
16                 content_id,
17                 transcription,
18                 os.path.basename(file_path),
19             )
20         )
21
22     for _ in tqdm(
23         as_completed(futures), total=len(futures), desc=" Processing audios"
24     ):
25         pass
26
27     print(f" Processed {len(metadata)} audio files.")

```

تعریف تابع

```

1 def process_audio_json(json_path: str, max_workers: int = 8):
2     """Process and store audio embeddings using metadata JSON."""

```

ورودی ها:

- *json_path*: مسیر فایل JSON، شامل متادیتا (اطلاعات مربوط به فایل های صوتی و متن آنها).
 - *max_workers*: تعداد نخ ها (*threads*) برای اجرای موازی (پیش فرض: ۸).
- خواندن فایل JSON، ورودی

```

1 with open(json_path, "r") as f:
2     metadata = json.load(f)

```

فایل JSON باز شده و محتویات آن در متغیر *metadata* ذخیره می شود.
ایجاد *ThreadPool* برای پردازش موازی

```

1 with ThreadPoolExecutor(max_workers=max_workers) as executor:

```

از *ThreadPoolExecutor* برای اجرای هم زمان چند پردازش استفاده می شود. تعداد نخ ها برابر مقدار *max_workers* تعیین می شود.
آماده سازی وظایف (Tasks)

```

1 futures = []
2 for item in metadata:
3     file_path = item.get("audio_path", "").strip()
4     transcription = item.get("transcription", "").strip()
5     if not (file_path and transcription):
6         continue

```

در این بخش: مسیر فایل صوتی (*audio_path*) و متن مربوطه (*transcription*) گرفته می شود. اگر مسیر یا متن خالی باشد، آن مورد نادیده گرفته می شود.
ارسال کارها به *ThreadPool*

```

1 content_id = str(item.get("id", os.path.basename(file_path)))
2 futures.append(
3     executor.submit(
4         store_audio_item,
5         content_id,
6         transcription,
7         os.path.basename(file_path),
8     )
9 )

```

برای هر فایل معتبر، تابع *store_audio_item* به صورت غیرهم زمان اجرا می شود.
ورودی های آن شامل موارد زیر است:

- *content_id*: شناسه ی فایل (از *id* یا نام فایل).

- *transcription*: متن گفتار استخراج شده.

- *os.path.basename(file_path)*: نام فایل بدون مسیر.

نمایش نوار پیشرفت


```

1 for _ in tqdm(
2     as_completed(futures), total=len(futures), desc=" Processing audios"
3 ):
4     pass

```

از *tqdm* برای نمایش نوار پیشرفت استفاده می‌شود. *as_completed(futures)* وضعیت تکمیل هر نخ را پیگیری می‌کند. خروجی نهایی در پایان، پیام موفقیت چاپ می‌شود: "Processed X audio files." که نشان می‌دهد چند فایل پردازش و ذخیره شده‌اند.

تابع *store_audio_item*

در تابع بالا از تابع زیر استفاده شده است. این تابع وظیفه دارد اطلاعات مربوط به هر فایل صوتی (شناسه، متن استخراج شده و *embedding*) را در پایگاه داده‌ی برداری *Weaviate* ذخیره نماید. در این فرآیند، بردار *embedding* از متن استخراج شده (و نه از خود فایل صوتی) تولید می‌شود.

```

1 def store_audio_item(item_id: str, transcription_text: str, audio_name: str = ""):
2     """Store audio transcription embedding and metadata in Weaviate."""
3     relative_path = f"/content/audio_dataset/{audio_name}"
4     embedding = get_embedding("text", transcription_text)
5     properties = {
6         "contentId": item_id,
7         "modality": "audio",
8         "filePath": relative_path,
9         "content": transcription_text,
10    }
11    collection.data.insert(properties=properties, vector=embedding.tolist())

```

تعریف تابع

```

1 def store_audio_item(item_id: str, transcription_text: str, audio_name: str = ""):
2     """Store audio transcription embedding and metadata in Weaviate."""

```

ورودی‌ها:

- *item_id*: شناسه‌ی منحصربه‌فرد فایل (می‌تواند از فایل *JSON* یا نام فایل گرفته شود).
- *transcription_text*: متنی که از گفتار صوتی استخراج شده است.

ساخت مسیر فایل (نسبی)

```

1 relative_path = f"/content/audio_dataset/{audio_name}"

```

تولید *embedding* از متن گفتار

```

1 embedding = get_embedding("text", transcription_text)

```

- تابع *get_embedding* یک بردار عددی (Embedding) از محتوای معنایی متن تولید می‌کند.
- پارامتر "text" نشان می‌دهد که نوع داده‌ی ورودی، متن است (نه تصویر یا صدا).
- این بردار برای متن استخراج شده تولید می‌شود، نه خود فایل صوتی.
- بردار *embedding* در مراحل بعدی برای جستجوی معنایی (*semantic search*) در *Weaviate* مورد استفاده قرار می‌گیرد.

آماده‌سازی متادیتا برای ذخیره

```

1 properties = {
2     "contentId": item_id,
3     "modality": "audio",
4     "filePath": relative_path,
5     "content": transcription_text,
6 }

```

- یک دیکشنری از ویژگی‌های شیء ساخته می‌شود که شامل اطلاعات توصیفی فایل است:
 - "contentId": شناسه‌ی منحصربه‌فرد فایل.
 - "modality": نوع داده

– "filePath": مسیر فایل صوتی در دیتاست.

– "content": متن استخراج شده از گفتار.

درج داده در پایگاه داده‌ی Weaviate

```
collection.data.insert(properties=properties, vector=embedding.tolist())
```

- داده‌ها در یک *Collection* (برای مثال "*Multimodal_Collection*") در *Weaviate* ذخیره می‌شوند.

- پارامترها شامل موارد زیر هستند:

- *properties*: متادیتا و اطلاعات توصیفی شامل شناسه، مسیر، نوع داده و متن.

- *vector*: بردار عددی *embedding* که به صورت لیست ذخیره می‌شود.

- به این ترتیب، هر داده‌ی صوتی در *Weaviate* شامل دو بخش اصلی است:

- متادیتا (ویژگی‌های توصیفی)

- بردار معنایی (*embedding*)

- این ساختار باعث می‌شود سیستم بتواند در مراحل بعدی بر اساس معنا و شباهت محتوایی میان فایل‌های صوتی، عملیات جست‌وجو و بازیابی مؤثر انجام دهد.

تابع *get_embedding*

```
1 def get_embedding(modality: str, input_data: Union[str, None]) -> np.ndarray:
2     mod = modality.lower()
3     if mod == "text":
4         if not isinstance(input_data, str):
5             raise ValueError("`input_data` must be a string.")
6
7         tokens = tokenizer([input_data]).to(DEVICE)
8         with torch.no_grad():
9             features = clip_model.encode_text(tokens)
10            features = features / features.norm(dim=-1, keepdim=True)
11            return features.detach().cpu().numpy().reshape(-1)
12     else:
13         raise ValueError("`modality` must be 'text'.")
```

در تابع *store_audio_item* از تابع زیر برای به دست آوردن *embedding* استفاده شده است. این تابع یک بردار *embedding* نرمال‌شده بر اساس نوع داده تولید می‌کند.

تعریف تابع و پارامترها

```
1 def get_embedding(modality: str, input_data: Union[str, None]) -> np.ndarray:
```

- *modality*: نوع داده

- *input_data*: داده‌ی ورودی برای تولید *embedding* (رشته‌ی متنی).

- خروجی: آرایه‌ی *NumPy* (*np.ndarray*) که بردار عددی متن را باز می‌گرداند.

توکنایز کردن متن

```
1 tokens = tokenizer([input_data]).to(DEVICE)
```

- متن با *tokenizer* مدل *CLIP* به توکن تبدیل می‌شود.

- توکن‌ها به دستگاه مناسب (*CPU* یا *GPU*) منتقل می‌شوند.

تولید *embedding* با مدل *CLIP*

```
1 with torch.no_grad():
2     features = clip_model.encode_text(tokens)
3     features = features / features.norm(dim=-1, keepdim=True)
```

- *torch.no_grad()*: جلوگیری از ذخیره‌سازی گرادیان‌ها برای صرفه‌جویی در حافظه و زمان.

- `clip_model.encode_text(tokens)`: تولید بردار عددی متن (*text embedding*) با مدل *CLIP*.
- `features / features.norm(dim=-1, keepdim=True)`: نرمال‌سازی بردار به طول یک.

بازگرداندن خروجی به *NumPy*

```
1 return features.detach().cpu().numpy().reshape(-1)
```

- `detach()`: جدا کردن بردار از گراف محاسباتی *PyTorch*.
- `cpu()`: انتقال داده به حافظه *CPU*.
- `numpy()`: تبدیل به آرایه‌ی *NumPy*.
- `reshape(-1)`: تبدیل به بردار یک‌بعدی.

۳.۶ فرآیند وارد کردن داده‌های متنی

به منظور وارد کردن داده‌های متنی به دیتابیس، لازم است متون مورد نظر در یک فایل با فرمت *CSV* ذخیره شده و در پوشه مربوطه قرار گیرد. در ادامه، مراحل اجرای اسکریپت و توضیح جزئیات کد ارائه شده است.

۱.۳.۶ اجرای اسکریپت

برای اجرای اسکریپت پایتون، ابتدا یک خط فرمان (*CMD*) را باز کرده و به مسیر ریشه‌ی پروژه بروید. سپس، با استفاده از دستور زیر، اسکریپت پایتون با نام *texts.py* را اجرا کنید:

```
1 python import_data_into_database/import_local_data/texts.py
```

۲.۳.۶ توضیحات جزئیات کد

این اسکریپت، وظیفه‌ی اتصال به دیتابیس، بارگذاری مدل *CLIP*، تعریف زنجیره‌ی ترجمه، پردازش داده‌ها (شامل ترجمه و تولید *embedding*) و ذخیره‌ی نهایی در دیتابیس *Weaviate* را بر عهده دارد.

۱. اتصال به دیتابیس *Weaviate* این بخش از کد، یک اتصال به سرویس محلی *Weaviate* که با *Docker* روی *localhost:8080* اجرا شده است، برقرار می‌سازد و سپس شیء *collection* مورد نظر را بازایی می‌کند.

```
1 try:
2     weaviate_client = weaviate.connect_to_local()
3     print(" Connected to Weaviate")
4     collection = weaviate_client.collections.get("Multimodal_Collection")
5 except Exception as e:
6     print(f" Weaviate connection failed: {e}")
7     weaviate_client = None
```

۲. تنظیمات مدل *CLIP* برای تولید *Embedding* کد زیر تنظیمات مربوط به مدل *CLIP* که برای تولید بردارهای عددی (*embedding*) داده‌های متنی و تصویری به کار می‌رود را انجام می‌دهد.

```
1 MODEL_NAME = "ViT-B-32"
2 PRETRAINED_LOCAL_PATH = (
3     "./open_clip_weights/ViT-B-32-openai/open_clip_model.safetensors"
4 )
5 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
6 clip_model, _, preprocess = open_clip.create_model_and_transforms(
7     MODEL_NAME, pretrained=PRETRAINED_LOCAL_PATH
8 )
9 tokenizer = open_clip.get_tokenizer(MODEL_NAME)
10 clip_model.to(DEVICE).eval()
```

توضیح پارامترها:

- `MODEL_NAME = "ViT-B-32"`: نام معماری مدل *CLIP* (مدل *Vision Transformer Base* با اندازه‌ی *patch* برابر با 32×32) مشخص می‌شود.
- `PRETRAINED_LOCAL_PATH`: مسیر فایل وزن‌های از پیش آموزش داده شده‌ی مدل را تعیین می‌کند.

- `DEVICE = "cuda" if torch.cuda.is_available() else "cpu"` در دسترس بودن *GPU* (*CUDA*) بررسی می‌شود. اگر *GPU* در دسترس باشد از آن، در غیر این صورت از *CPU* استفاده خواهد شد.
- `open_clip.create_model_and_transforms`: مدل *CLIP* را از مسیر داده‌شده بارگذاری کرده و تابع پیش‌پردازش تصاویر (*preprocess*) را برمی‌گرداند.
- `open_clip.get_tokenizer`: توکنایزر مخصوص همان مدل برای تبدیل متون به توکن‌های عددی قابل پردازش توسط مدل را بازپایی می‌کند.
- `clip_model.to(DEVICE).eval()`: مدل را روی دستگاه تعیین شده منتقل کرده و در حالت ارزیابی (*eval*) قرار می‌دهد تا برای پیش‌بینی آماده شود.

۳. بارگذاری کلید *API* برای ترجمه از آنجا که داده‌های متنی فارسی باید توسط یک مدل زبانی بزرگ (*LLM*) ترجمه شوند، نیاز است کلید *API* دریافتی از *OpenRouter* که در فایل `env` ذخیره شده است، بارگذاری شود.

```
1 load_dotenv()
2 API_KEY = os.getenv("LLM_API_KEY")
3 if not API_KEY:
4     raise ValueError(
5         "Please set your API key in the environment variable OPENROUTER_API_KEY."
6     )
```

توضیح:

- `load_dotenv()`: این دستور فایل `env` را خوانده و تمام متغیرهای محیطی (*environment variables*) داخل آن را در محیط اجرایی پایتون بارگذاری می‌کند.
- `API_KEY = os.getenv("LLM_API_KEY")`: با استفاده از تابع `os.getenv()`، مقدار متغیر محیطی با نام `"LLM_API_KEY"` خوانده و در متغیر `API_KEY` ذخیره می‌شود.

۴. تعریف زنجیره‌ی پردازش (*Pipeline*) برای ترجمه کد زیر یک زنجیره‌ی پردازش هوشمند برای ترجمه‌ی خودکار متن فارسی به انگلیسی ایجاد می‌کند که شامل مراحل دریافت متن فارسی، آماده‌سازی *prompt* ترجمه، ارسال به مدل *GPT* از طریق *OpenRouter* و دریافت ترجمه‌ی خالص است.

```
1 LLM_model = ChatOpenAI(
2     model="openai/gpt-oss-20b:free",
3     temperature=0.1,
4     openai_api_base="https://openrouter.ai/api/v1",
5     openai_api_key=API_KEY,
6 )
7 prompt = ChatPromptTemplate.from_messages(
8     [
9         (
10            "system",
11            "Translate the following Persian (Farsi) text to clear, natural English.\n"
12            "Return only the translation (no extra explanation).\n\n"
13            "Persian text:\n{text}",
14        )
15    ]
16 )
17 parser = StrOutputParser()
18 vision_chain = prompt | LLM_model | parser
```

تشریح اجزای زنجیره‌ی ترجمه:

۱. ساخت مدل زبانی (*LLM*)

```
1 LLM_model = ChatOpenAI(
2     model="openai/gpt-oss-20b:free",
3     temperature=0.1,
4     openai_api_base="https://openrouter.ai/api/v1",
5     openai_api_key=API_KEY,
6 )
7
```

- کارکرد: ایجاد یک شیء مدل زبانی (*LLM*) از کلاس *ChatOpenAI* (کتابخانه‌ی *LangChain*) برای برقراری ارتباط با مدل *GPT* از طریق *API*.

- پارامترها:

- (آ) `model="openai/gpt-oss-20b:free"`: نام مدل زبانی که از `OpenRouter` استفاده می‌کند
- (ب) `temperature=0.1`: درجه‌ی تصادفی بودن پاسخ‌ها را مشخص می‌کند (هرچه کمتر، پاسخ‌ها دقیق‌تر و تکراری‌تر).
- (ج) `openai_api_base="https://openrouter.ai/api/v1"`: آدرس پایه‌ی `API` سرویس `OpenRouter` (جایگزین مستقیم `OpenAI API`).
- (د) `openai_api_key=API_KEY`: کلید `API` شما برای احراز هویت در `OpenRouter`.

۲. تعریف قالب پیام (prompt template)

```
1 prompt = ChatPromptTemplate.from_messages(
2     [
3         (
4             "system",
5             "Translate the following Persian (Farsi) text to clear, natural
6             English.\n"
7             "Return only the translation (no extra explanation).\n\n"
8             "Persian text:\n{text}",
9         )
10    ]
11 )
```

- کارکرد: ساخت یک الگو برای پیام ورودی که به مدل داده می‌شود.
- توضیح جزئی:
 - (آ) `"system"`: نوع پیام (به مدل دستور می‌دهد چگونه رفتار کند).
 - (ب) پیام داخل کوتیشن: دستور ترجمه است؛ از مدل خواسته شده: متن فارسی را ترجمه کند، فقط ترجمه را برگرداند (بدون توضیح اضافه).
 - (ج) `{text}`: متغیری است که هنگام اجرا با متن واقعی جایگزین می‌شود.

۳. ساخت مفسر خروجی (parser)

```
1 parser = StrOutputParser()
2
```

- کارکرد: این مفسر خروجی خام مدل را گرفته و به رشته‌ی متنی نهایی (`string`) تبدیل می‌کند. به عبارتی، خروجی `JSON` یا ساختاریافته‌ی مدل را ساده کرده و فقط متن ترجمه را برمی‌گرداند.

۴. ساخت زنجیره‌ی پردازش (chain)

```
1 vision_chain = prompt | LLM_model | parser
2
```

- کارکرد: با استفاده از اپراتور `pipe` (`|`) در `LangChain`، سه جزء را به صورت مرحله‌به‌مرحله به هم وصل می‌کند: `prompt` (قالب ورودی)، `LLM_model` (مدل زبانی)، و `parser` (مفسر خروجی).
- نتیجه: `vision_chain` یک زنجیره‌ی کامل است که با دریافت یک متن فارسی، ترجمه‌ی انگلیسی نهایی را برمی‌گرداند.

۵. بارگذاری و پردازش داده‌های متنی بخش اصلی اسکریپت که فایل ورودی `CSV` را مشخص کرده و فرآیند دو مرحله‌ای پردازش را آغاز می‌کند:

```
1 if __name__ == "__main__":
2     # Step 1: Input file (contains Persian + English)
3     input_csv = "./content/text_dataset/text_data.csv"
4
5     # Step 2: Prepare final English-only CSV
6     final_csv = prepare_final_csv(input_csv, text_column="text", workers=2)
7
8     # Step 3: Generate and store embeddings in Weaviate
9     process_texts(final_csv)
10    weaviate_client.close()
```

فرآیند پردازش داده‌های متنی در دو مرحله انجام می‌شود:

۱. مرحله اول: ترجمه متون به انگلیسی: تمام داده‌های متنی (فارسی و انگلیسی) پردازش شده و تمام داده‌ها به انگلیسی ترجمه می‌شوند.

۲. مرحله دوم: تولید *embedding* و ذخیره در دیتابیس: برای داده‌های انگلیسی تولید شده، *embedding* محاسبه شده و در دیتابیس *Weaviate* ذخیره می‌شود.

مرحله اول: ترجمه متون به انگلیسی با استفاده از تابع *prepare_final_csv* تمام داده‌ها را به انگلیسی ترجمه می‌کنیم. اگر متن انگلیسی باشد نیاز به ترجمه نیست، اما اگر فارسی باشد ابتدا با استفاده از *LLM* آن را ترجمه می‌کنیم. تابع *prepare_final_csv()* سه ورودی (پارامتر) می‌گیرد:

۱. *input_csv*: مسیر فایل *CSV* ورودی (حاوی متن‌ها، فارسی و انگلیسی).

۲. *text_column="text"*: نام ستونی از *CSV* که متن‌ها در آن قرار دارند.

۳. *workers=2*: تعداد (*threads*) یا پردازنده‌های موازی برای سرعت‌دهی به پردازش.

خروجی این تابع یک فایل *CSV* جدید خواهد بود که در آن تمام داده‌ها به انگلیسی ترجمه شده‌اند.

توضیح تابع *prepare_final_csv* این تابع با استفاده از *LLM* متون فارسی را ترجمه می‌کند و در نهایت یک فایل *CSV* جدید حاوی تمام متون ترجمه شده به زبان انگلیسی تولید می‌کند.

```
1 def prepare_final_csv(  
2     input_csv: str, text_column="text", translated_column="translated", workers=2  
3 ):  
4     """  
5     Process a CSV file where each row can be Persian or English.  
6     - English rows are added directly to final.csv  
7     - Persian rows are translated and also added to the same 'text' column  
8     """  
9     if not os.path.isfile(input_csv):  
10         raise FileNotFoundError(f"File '{input_csv}' does not exist.")  
11  
12     df = pd.read_csv(input_csv)  
13  
14     if text_column not in df.columns:  
15         raise ValueError(  
16             f"Column '{text_column}' does not exist in the file. Columns:  
17             {df.columns.tolist()}"  
18         )  
19  
20     # Split English and Persian rows  
21     english_rows = df[~df[text_column].apply(is_persian)]  
22     persian_rows = df[df[text_column].apply(is_persian)]  
23  
24     # Paths  
25     folder, filename = os.path.split(input_csv)  
26     name, ext = os.path.splitext(filename)  
27     final_csv = os.path.join(folder, f"{name}_final.csv")  
28     temp_csv = os.path.join(folder, f"{name}_temp.csv")  
29  
30     # Save English rows directly to final.csv  
31     english_rows[[text_column]].to_csv(final_csv, index=False, encoding="utf-8-sig")  
32  
33     if persian_rows.empty:  
34         print("No Persian text found. Final file saved.")  
35         return final_csv  
36  
37     # Save Persian rows to temp.csv for translation  
38     persian_rows[[text_column]].to_csv(temp_csv, index=False, encoding="utf-8-sig")  
39  
40     # Translate Persian rows  
41     translated_csv = translate_csv(temp_csv, text_column=text_column, workers=workers)  
42  
43     # Load translated text  
44     translated_df = pd.read_csv(translated_csv)  
45  
46     # We only need the translated text under the same column name 'text'
```

```

46 translated_text = translated_df[[translated_column]].rename(
47     columns={translated_column: text_column}
48 )
49
50 # Combine English + translated rows (all under one column 'text')
51 final_df = pd.concat(
52     [english_rows[[text_column]], translated_text], ignore_index=True
53 )
54
55 # Save final CSV (only one column: 'text')
56 final_df.to_csv(final_csv, index=False, encoding="utf-8-sig")
57
58 # Clean up
59 os.remove(temp_csv)
60 os.remove(translated_csv)
61
62 print(f" Final CSV prepared (single 'text' column): {final_csv}")
63 return final_csv

```

توضیح کد `prepare_final_csv`

۱. تعریف تابع و پارامترها:

```

1 def prepare_final_csv(
2     input_csv: str, text_column="text", translated_column="translated", workers=2
3 ):
4

```

- `input_csv`: مسیر فایل *CSV* ورودی.
- `text_column`: نام ستونی که متن‌ها در آن هستند (پیش‌فرض `"text"`).
- `translated_column`: نام ستونی که متن ترجمه‌شده بعداً در آن ذخیره می‌شود (پیش‌فرض `"translated"`).
- `workers`: تعداد رشته‌های موازی برای ترجمه (پیش‌فرض ۲).

۲. بررسی وجود فایل:

```

1 if not os.path.isfile(input_csv):
2     raise FileNotFoundError(f"File '{input_csv}' does not exist.")
3

```

- اگر فایل ورودی وجود نداشته باشد، خطا می‌دهد و برنامه متوقف می‌شود.

۳. خواندن فایل *CSV*:

```

1 df = pd.read_csv(input_csv)
2

```

- فایل *CSV* را با *pandas* می‌خواند و آن را در *DataFrame* ذخیره می‌کند.

۴. بررسی وجود ستون متن:

```

1 if text_column not in df.columns:
2     raise ValueError(
3         f"Column '{text_column}' does not exist in the file. Columns:
4         {df.columns.tolist()}"
5     )

```

- بررسی می‌کند که آیا ستونی با نام `"text"` در فایل وجود دارد یا نه؛ اگر نه، خطا می‌دهد.

۵. جدا کردن سطرهای فارسی و انگلیسی:

```

1 english_rows = df[~df[text_column].apply(is_persian)]
2 persian_rows = df[df[text_column].apply(is_persian)]
3

```

- تابع `is_persian()` برای هر سطر بررسی می‌کند که آیا متن فارسی است یا نه.
- سطرهای انگلیسی در `english_rows` و سطرهای فارسی در `persian_rows` ذخیره می‌شوند.

۶. تعریف مسیر فایل‌های خروجی:

```
1 folder, filename = os.path.split(input_csv)
2 name, ext = os.path.splitext(filename)
3 final_csv = os.path.join(folder, f"{name}_final.csv")
4 temp_csv = os.path.join(folder, f"{name}_temp.csv")
5
```

- مسیر فایل ورودی را تجزیه می‌کند و مسیر فایل‌های موقت و نهایی را می‌سازد: `*_temp.csv` (شامل فقط سطرهای فارسی برای ترجمه) و `*_final.csv` (شامل خروجی نهایی انگلیسی).

۷. ذخیره‌ی مستقیم سطرهای انگلیسی:

```
1 english_rows[[text_column]].to_csv(final_csv, index=False, encoding="utf-8-sig")
2
```

- سطرهای انگلیسی بدون تغییر مستقیماً در فایل نهایی ذخیره می‌شوند.

۸. در صورت نبود متن فارسی:

```
1 if persian_rows.empty:
2     print("No Persian text found. Final file saved.")
3     return final_csv
4
```

- اگر در داده‌ها هیچ متن فارسی نباشد، مستقیماً فایل نهایی را برمی‌گرداند.

۹. ذخیره‌ی سطرهای فارسی برای ترجمه:

```
1 persian_rows[[text_column]].to_csv(temp_csv, index=False, encoding="utf-8-sig")
2
```

- تمام سطرهای فارسی را در یک فایل موقت (`temp_csv`) ذخیره می‌کند تا برای ترجمه استفاده شود.

۱۰. ترجمه‌ی متون فارسی:

```
1 translated_csv = translate_csv(temp_csv, text_column=text_column, workers=workers)
2
```

- فایل موقت فارسی را با تابع `translate_csv()` ترجمه می‌کند. این تابع به‌صورت موازی (`workers=2`) اجرا می‌شود و خروجی را به‌صورت `CSV` برمی‌گرداند.

۱۱. خواندن خروجی ترجمه:

```
1 translated_df = pd.read_csv(translated_csv)
2
```

- فایل ترجمه‌شده را می‌خواند و در `DataFrame` جدید ذخیره می‌کند.

۱۲. تغییر نام ستون ترجمه به `"text"`:

```
1 translated_text = translated_df[[translated_column]].rename(
2     columns={translated_column: text_column}
3 )
4
```

- برای یک‌دست شدن، ستون `"translated"` را به `"text"` تغییر نام می‌دهد تا با سطرهای انگلیسی هماهنگ شود.

۱۳. ترکیب داده‌های انگلیسی و ترجمه‌شده:


```

1 final_df = pd.concat(
2     [english_rows[[text_column]], translated_text], ignore_index=True
3 )
4

```

• هر دو مجموعه (انگلیسی و ترجمه شده) را زیر هم در یک *DataFrame* ادغام می کند.

۱۴. ذخیره ی فایل نهایی:

```

1 final_df.to_csv(final_csv, index=False, encoding="utf-8-sig")
2

```

• همه ی متن ها (الان همه انگلیسی هستند) را در یک فایل *CSV* نهایی ذخیره می کند.

۱۵. حذف فایل های موقت:

```

1 os.remove(temp_csv)
2 os.remove(translated_csv)
3

```

• فایل های موقتی (فارسی و ترجمه شده) را برای تمیزی محیط حذف می کند.

۱۶. پیام نهایی و بازگرداندن مسیر:

```

1 print(f" Final CSV prepared (single 'text' column): {final_csv}")
2 return final_csv
3

```

• آدرس فایل نهایی را چاپ و برمی گرداند.

تابع *is_persian* برای تشخیص متون فارسی این تابع با استفاده از عبارت منظم (*Regex*) بررسی می کند که آیا حداقل یک کاراکتر فارسی در متن وجود دارد یا خیر.

```

1 def is_persian(text: str) -> bool:
2     """Return True if the text contains at least one Persian character."""
3     if not text or not isinstance(text, str):
4         return False
5     return bool(re.search(r"[\u0600-\u06FF]", text))

```

منطق: این تابع با استفاده از عبارت منظم (*re.search(r"[\u0600-\u06FF]", text)*) بررسی می کند که آیا حداقل یک کاراکتر در محدوده ی یونیکد کاراکترهای فارسی و عربی (0x0600 تا 0x06FF) در متن وجود دارد یا خیر. در صورت مثبت بودن، مقدار *True* بازگردانده می شود.

• از ماژول *re* (عبارت منظم) استفاده می کند.

• الگوی *r"[\u0600-\u06FF]"*:

۱. محدوده ی یونیکد کاراکترهای فارسی و عربی (از 0x0600 تا 0x06FF) را مشخص می کند.

۲. یعنی اگر حداقل یک کاراکتر در این محدوده وجود داشته باشد، تطبیق پیدا می کند.

• *re.search()*: اولین تطبیق را پیدا می کند و اگر باشد، شیء *match* برمی گرداند.

• *bool(...)*: تبدیل نتیجه به *True* یا *False*.

تابع `translate_csv` برای ترجمه‌ی موازی این تابع مسئول ترجمه‌ی متون یک ستون مشخص از فایل `CSV` به صورت موازی است و خروجی را در یک فایل جدید (`<original_name>_translated.csv`) ذخیره می‌کند.

```

1 def translate_csv(input_csv: str, text_column="text", workers=4):
2     """
3     Takes a CSV file, translates the specified text column,
4     and saves the output in the same folder as the input file with
5     the name <original_name>_translated.csv.
6     """
7     if not os.path.isfile(input_csv):
8         raise FileNotFoundError(f"File '{input_csv}' does not exist.")
9
10    df = pd.read_csv(input_csv)
11
12    if text_column not in df.columns:
13        raise ValueError(
14            f"Column '{text_column}' does not exist in the file. Columns:
15            {df.columns.tolist()}"
16        )
17
18    texts = df[text_column].fillna("").astype(str).tolist()
19    results = [""] * len(texts)
20
21    # Use ThreadPoolExecutor to speed up translation (optional)
22    if workers > 1:
23        with ThreadPoolExecutor(max_workers=workers) as ex:
24            futures = {
25                ex.submit(translate_text, texts[i]): i for i in range(len(texts))
26            }
27            for fut in tqdm(
28                as_completed(futures), total=len(futures), desc="Translating"
29            ):
30                idx = futures[fut]
31                try:
32                    results[idx] = fut.result()
33                except Exception as e:
34                    print(f"Error at index {idx}: {e}")
35                    results[idx] = ""
36            else:
37                for i, t in enumerate(tqdm(texts, desc="Translating")):
38                    results[i] = translate_text(t)
39
40    df["translated"] = results
41
42    # Save output file in the same folder as the input file
43    folder, filename = os.path.split(input_csv)
44    name, ext = os.path.splitext(filename)
45    output_csv = os.path.join(folder, f"{name}_translated{ext}")
46
47    df.to_csv(output_csv, index=False, encoding="utf-8-sig")
48    print(f"Translation complete. Saved to {output_csv}")
49    return output_csv

```

توضیح: این تابع از `ThreadPoolExecutor` برای اجرای موازی تابع `translate_text` بر روی متون استفاده می‌کند تا سرعت فرآیند ترجمه افزایش یابد. نوار پیشرفت (`tqdm`) نیز برای نمایش وضعیت ترجمه استفاده می‌شود.

توضیح کد `translate_csv`

۱. تعریف تابع و پارامترها:

```

1 def translate_csv(input_csv: str, text_column="text", workers=4):
2

```

- `input_csv`: مسیر فایل `CSV` ورودی.
- `text_column`: نام ستونی که متن‌ها در آن قرار دارند (پیش‌فرض `"text"`).
- `workers`: تعداد رشته‌های موازی (`threads`) برای ترجمه (پیش‌فرض ۴).

۲. بررسی وجود فایل:

```

1 if not os.path.isfile(input_csv):
2     raise FileNotFoundError(f"File '{input_csv}' does not exist.")
3

```

• اگر فایل ورودی وجود نداشته باشد، خطا می‌دهد.

۳. خواندن فایل CSV:

```

1 df = pd.read_csv(input_csv)
2

```

• فایل CSV را با *pandas* خوانده و در یک *DataFrame* ذخیره می‌کند.

۴. بررسی وجود ستون متن:

```

1 if text_column not in df.columns:
2     raise ValueError(
3         f"Column '{text_column}' does not exist in the file. Columns:
4         {df.columns.tolist()}"
5     )

```

• مطمئن می‌شود ستونی که قرار است ترجمه شود در فایل وجود داشته باشد؛ در غیر این صورت خطا می‌دهد.

۵. آماده‌سازی متن‌ها:

```

1 texts = df[text_column].fillna("").astype(str).tolist()
2 results = [""] * len(texts)
3

```

• *texts*: تمام مقادیر ستون متن را به رشته تبدیل کرده و جای مقادیر خالی (*NaN*) را با رشته خالی پر می‌کند.

• *results*: آرایه‌ای برای ذخیره‌ی نتایج ترجمه، هم‌اندازه با تعداد سطرها.

۶. ترجمه موازی (با *ThreadPoolExecutor*):

```

1 if workers > 1:
2     with ThreadPoolExecutor(max_workers=workers) as ex:
3         futures = {
4             ex.submit(translate_text, texts[i]): i for i in range(len(texts))
5         }
6         for fut in tqdm(
7             as_completed(futures), total=len(futures), desc="Translating"
8         ):
9             idx = futures[fut]
10            try:
11                results[idx] = fut.result()
12            except Exception as e:
13                print(f" Error at index {idx}: {e}")
14                results[idx] = ""
15

```

• اگر تعداد *workers* بیشتر از ۱ باشد، ترجمه به صورت موازی انجام می‌شود.

• هر متن به یک *thread* فرستاده می‌شود (*translate_text*).

• *tqdm* برای نمایش نوار پیشرفت استفاده شده است.

• در صورت بروز خطا، رشته خالی جایگزین می‌شود و پیام خطا چاپ می‌شود.

۷. ترجمه تک‌نخی (بدون موازی):

```

1 else:
2     for i, t in enumerate(tqdm(texts, desc="Translating")):
3         results[i] = translate_text(t)
4

```

- اگر `workers=1` باشد، ترجمه به صورت تک‌نخی انجام می‌شود، ولی باز هم نوار پیشرفت نشان داده می‌شود.

۸. اضافه کردن ستون ترجمه:

```
1 df["translated"] = results
2
```

- ستون جدید `"translated"` به `DataFrame` اضافه می‌شود و نتایج ترجمه در آن قرار می‌گیرند.

۹. تعیین مسیر فایل خروجی:

```
1 folder, filename = os.path.split(input_csv)
2 name, ext = os.path.splitext(filename)
3 output_csv = os.path.join(folder, f"{name}_translated{ext}")
4
```

- مسیر پوشه و نام فایل ورودی استخراج می‌شود.
- مسیر فایل خروجی به صورت `<original_name>_translated.csv` ساخته می‌شود.

۱۰. ذخیره `CSV` ترجمه شده:

```
1 df.to_csv(output_csv, index=False, encoding="utf-8-sig")
2 print(f" Translation complete. Saved to {output_csv}")
3
```

- `DataFrame` شامل متن اصلی و ترجمه در فایل جدید ذخیره می‌شود.
- پیام موفقیت چاپ می‌شود.

۱۱. بازگرداندن مسیر فایل خروجی:

```
1 return output_csv
2
```

- مسیر فایل `CSV` ترجمه شده را برمی‌گرداند.

تابع `translate_text` برای ترجمه‌ی یک آیت‌متنی این تابع مسئول ترجمه‌ی یک رشته‌ی متنی با استفاده از زنجیره‌ی `vision_chain` و مدیریت خطاهای احتمالی از طریق تلاش مجدد (`retries`) است.

```
1 def translate_text(text: str, retries=2, pause=1.0) -> str:
2     # ----- Translate a single text -----
3     if not isinstance(text, str) or not text.strip():
4         return ""
5
6     prompt_input = {"text": text}
7
8     attempt = 0
9     while attempt <= retries:
10         try:
11             translated = vision_chain.invoke(prompt_input)
12             return translated.strip()
13         except Exception as e:
14             attempt += 1
15             print(f" Translation error (attempt {attempt}): {e}")
16             if attempt > retries:
17                 return ""
18             time.sleep(pause * attempt)
```

توضیح کد `translate_text`

۱. تعریف تابع و پارامترها:

```
1 def translate_text(text: str, retries=2, pause=1.0) -> str:
```

2

- `text`: رشته‌ای که قرار است ترجمه شود.
- `retries`: تعداد دفعات تلاش مجدد در صورت خطا (پیش‌فرض ۲).
- `pause`: مدت زمان مکث بین تلاش‌ها بر حسب ثانیه (پیش‌فرض ۱ ثانیه).
- خروجی: رشته‌ی ترجمه‌شده (`str`).

۲. آماده‌سازی ورودی برای زنجیره ترجمه:

```
1 prompt_input = {"text": text}
2 translated = vision_chain.invoke(prompt_input)
3 return translated.strip()
```

4

- متن ورودی در یک دیکشنری با کلید `"text"` قرار داده می‌شود.
- `vision_chain.invoke(prompt_input)`: متن را به زنجیره‌ی ترجمه می‌دهد و ترجمه نهایی را می‌گیرد.
- در صورت موفقیت، متن ترجمه‌شده را با `strip()` برمی‌گرداند (حذف فاصله‌های اضافی).

۳. مدیریت خطا و تلاش مجدد:

```
1 except Exception as e:
2     attempt += 1
3     print(f" Translation error (attempt {attempt}): {e}")
4     if attempt > retries:
5         return ""
6     time.sleep(pause * attempt)
```

7

- اگر خطایی رخ دهد (مانند قطع اتصال `API` یا خطای مدل):
 - (آ) تعداد تلاش‌ها افزایش می‌یابد (`attempt += 1`).
 - (ب) پیام خطا چاپ می‌شود.
 - (ج) اگر تعداد تلاش‌ها بیشتر از مقدار `retries` باشد → رشته خالی برگردانده می‌شود.
 - (د) قبل از تلاش مجدد، تابع به مدت `attempt * pause` ثانیه مکث می‌کند (استراتژی افزایش فاصله بین تلاش‌ها).

خلاصه‌ی مرحله اول: حال مرحله ۱ به پایان می‌رسد و در پایان این مرحله ما دارای یک فایل `CSV` جدید هستیم که تمام داده‌های آن به انگلیسی ترجمه شده‌اند. این فایل در همان پوشه‌ای قرار دارد که فایل `CSV` اصلی قرار داشت.

مرحله دوم: تولید `Embedding` و ذخیره در `Weaviate` پس از اتمام مرحله اول، فایل `CSV` نهایی حاوی متون انگلیسی آماده شده است. خط زیر مسئول شروع مرحله دوم است:

```
1 process_texts(final_csv)
```

2

این دستور تابع `process_texts` را با مسیر فایل `CSV` نهایی (`final_csv`) فراخوانی می‌کند.

تابع `process_texts` این تابع مسئول اجرای موازی تولید `embedding` و ذخیره‌ی آیتم‌های متنی در دیتابیس است.

```
1 def process_texts(input_csv: str, max_workers: int = 8):
2     """Process and store text embeddings in parallel."""
3     df = pd.read_csv(input_csv)
4     texts = df["text"].astype(str).tolist()
5
6     with ThreadPoolExecutor(max_workers=max_workers) as executor:
7         futures = [
8             executor.submit(store_text_item, f"text_{i+1}", text)
9             for i, text in enumerate(texts)
10        ]
```

```

11     for _ in tqdm(
12         as_completed(futures), total=len(futures), desc=" Processing texts"
13     ):
14         pass
15
16     print(f" Processed {len(texts)} text items.")

```

توضیح کد *process_texts*

۱. تعریف تابع و پارامترها:

```

1 def process_texts(input_csv: str, max_workers: int = 8):
2

```

- *input_csv*: مسیر فایل *CSV* ورودی که شامل ستون *"text"* است.
- *max_workers*: حداکثر تعداد رشته‌های موازی (*threads*) برای پردازش متون (پیش‌فرض ۸).

۲. خواندن فایل *CSV* و آماده‌سازی متن‌ها:

```

1 df = pd.read_csv(input_csv)
2 texts = df["text"].astype(str).tolist()
3

```

- فایل *CSV* ورودی را با *pandas* می‌خواند.
- ستون *"text"* را به رشته تبدیل و به لیست متن‌ها (*texts*) تبدیل می‌کند.

۳. اجرای پردازش موازی:

```

1 with ThreadPoolExecutor(max_workers=max_workers) as executor:
2     futures = [
3         executor.submit(store_text_item, f"text_{i+1}", text)
4         for i, text in enumerate(texts)
5     ]
6

```

- از *ThreadPoolExecutor* برای اجرای همزمان چند عملیات استفاده می‌شود.
- برای هر متن، تابع *store_text_item()* فراخوانی می‌شود.
- به هر متن یک شناسه یکتا (*text_1, text_2, ...*) داده می‌شود.
- هر فراخوانی به *executor.submit()* داده می‌شود تا در یک رشته‌ی مجزا اجرا شود.
- نتیجه‌ها در *futures* ذخیره می‌شوند (لیست *future objects*).

۴. نمایش پیشرفت:

```

1 for _ in tqdm(
2     as_completed(futures), total=len(futures), desc=" Processing texts"
3 ):
4     pass
5

```

- حلقه بر روی *as_completed(futures)* می‌چرخد تا منتظر بماند تمام وظایف موازی تکمیل شوند.
- *tqdm* برای نمایش نوار پیشرفت استفاده شده است.

تابع *store_text_item* این تابع عملیات تولید *embedding* برای یک متن و ذخیره‌ی آن همراه با متادیتا در دیتابیس *Weaviate* را انجام می‌دهد.

```

1 def store_text_item(item_id: str, text_data: str):
2     """Store text embedding and metadata in Weaviate."""
3     embedding = get_embedding("text", text_data)
4     properties = {
5         "contentId": item_id,
6         "modality": "text",
7         "filePath": "",
8         "content": text_data,
9     }
10    collection.data.insert(properties=properties, vector=embedding.tolist())

```

توضیح کد `store_text_item`

- تعریف تابع و پارامترها:

- `item_id`: شناسه یکتا برای متن (مثلاً `"text_1"`).

- `text_data`: خود متن مورد نظر برای پردازش.

- تولید `embedding`:

- `embedding = get_embedding("text", text_data)`: بردار عددی (`embedding`) متن را تولید می‌کند.

- این بردار، نمایش معنایی متن در فضای چندبعدی است و برای جستجوی معنایی و بازیابی متون استفاده می‌شود.

- آماده‌سازی متادیتا:

- دیکشنری `properties` شامل اطلاعات جانبی است که همراه `embedding` ذخیره می‌شود: `contentId` (شناسه یکتا)، `modality` (نوع داده)، `filePath` (مسیر فایل منبع، که در اینجا خالی است) و `content` (خود متن).

- ذخیره در `Weaviate`:

- `collection.data.insert()`: داده را در `collection` مشخص در `Weaviate` ذخیره می‌کند.

- پارامتر `properties`: متادیتای متن.

- پارامتر `vector`: بردار `embedding` متن (که با `tolist()` به لیست تبدیل می‌شود تا قابل ذخیره باشد).

۵- تابع `get_embedding` این تابع مسئول تولید بردار `embedding` نرمالایز شده توسط مدل `CLIP` برای یک داده‌ی متنی است.

```
1 def get_embedding(modality: str, input_data: Union[str, None]) -> np.ndarray:
2     """Wrapper for modality-specific embedding."""
3     mod = modality.lower()
4     if mod == "text":
5         """Return normalized CLIP text embedding."""
6         if not isinstance(input_data, str):
7             raise ValueError("`text` must be a string.")
8
9         tokens = tokenizer([input_data]).to(DEVICE)
10        with torch.no_grad():
11            features = clip_model.encode_text(tokens)
12            features = features / features.norm(dim=-1, keepdim=True)
13        return features.detach().cpu().numpy().reshape(-1)
14    else:
15        raise ValueError("`modality` must be 'text'.")
```

توضیح کد `get_embedding`

- تعریف تابع و پارامترها:

- `modality`: نوع داده.

- `input_data`: داده ورودی برای تولید `embedding`.

- خروجی: آرایه‌ی `NumPy` (`np.ndarray`) که بردار عددی متن را باز می‌گرداند.

- توکنایز کردن متن:

- `tokens = tokenizer([input_data]).to(DEVICE)`: متن را با `tokenizer` مدل `CLIP` به توکن تبدیل کرده و به دستگاه مناسب (`cpu` یا `cuda`) منتقل می‌کند.

- تولید `embedding` با مدل `CLIP`:

- `with torch.no_grad()`: جلوگیری از ذخیره‌سازی گرادینان‌ها برای صرفه‌جویی در حافظه و زمان.

- `features = clip_model.encode_text(tokens)`: تولید بردار عددی متن (`text embedding`) با مدل `CLIP`.

- `features = features / features.norm(dim=-1, keepdim=True)`: نرمال‌سازی بردار (طول بردار برابر ۱ می‌شود).

- بازگرداندن خروجی به `NumPy`:

- `features.detach().cpu().numpy().reshape(-1)`

۱. `detach()`: جدا کردن از گراف محاسباتی `PyTorch`.
۲. `cpu()`: انتقال به حافظه `CPU`.
۳. `numpy()`: تبدیل به آرایه `NumPy`.
۴. `reshape(-1)`: تبدیل به بردار یک‌بعدی.

۷ بخش بک اند اصلی

کد اصلی بک اند داخل پوشه `server` قرار دارد در زیربخش های بعدی به توضیح فایل های موجود رد پوشه `server` و عملکرد آنها پرداخته خواهد شد

۱.۷ فایل `config.py`

قبل از آغاز عملیات اجرایی، لازم است فایل `config.py` واقع در پوشه `server` بررسی و در صورت لزوم تنظیم گردد. این فایل شامل مجموعه ای از تنظیمات اولیه و پیکربندی ها برای تعیین مسیر فایل ها و مشخص نمودن انواع مدل های مورد استفاده در پروژه است. در صورتی که روال های تعیین شده در بخش های پیشین پروژه به دقت دنبال شده باشد، نیازی به اعمال تغییرات در این مسیرها و تنظیمات نخواهد بود. در غیر این صورت، موارد درخواستی باید توسط کاربر تنظیم و به روزرسانی شوند. کد زیر، محتوای کامل فایل `config.py` را نشان می دهد:

```
1 import os
2 import torch
3 from dotenv import load_dotenv
4 from pathlib import Path
5
6 PROJECT_ROOT_PATH = Path.cwd().as_posix()
7 print(f"Project root path: {PROJECT_ROOT_PATH}")
8
9 WEAVIATE_COLLECTION_NAME = "Multimodal_Collection"
10
11 # === Load environment variables ===
12 load_dotenv() # Looks for .env in current working directory or parents
13
14 # === Device setup ===
15 DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
16
17 # === CLIP Model Config ===
18 CLIP_MODEL_NAME = "ViT-B-32"
19 PRETRAINED_LOCAL_PATH =
20     f"{PROJECT_ROOT_PATH}/open_clip_weights/ViT-B-32-openai/open_clip_model.safetensors"
21
22 # === Whisper Model Config ===
23 WHISPER_MODEL = "small"
24
25 # === LLM Config ===
26 LLM_MODEL_NAME = "openai/gpt-5-image-mini"
27 LLM_API_BASE = "https://openrouter.ai/api/v1"
28 LLM_API_KEY = os.getenv("LLM_API_KEY") # <- Loaded from .env
29
30 if not LLM_API_KEY:
31     raise ValueError("Missing LLM_API_KEY in .env file!")
32
33 print("Environment key loaded successfully!")
```

این قطعه کد به منظور پیکربندی (*Configuration*) یک پروژه ی چندحالتی (*Multimodal*) طراحی شده است که بر مبنای مدل های زبان بزرگ (*LLM*) و مدل های پردازش تصویر (مانند *Multimodal RAG* یا چت بات چندحالتی) عمل می کند.

۱.۱.۷ مسیر ریشه پروژه

```
1 PROJECT_ROOT_PATH = Path.cwd().as_posix()
2 print(f"Project root path: {PROJECT_ROOT_PATH}")
```

- `Path.cwd()`: مسیر جاری عملیاتی (*Current Working Directory*) را بازایی می نماید.
- `as_posix()`: مسیر را به فرم استاندارد یونیکس (/) تبدیل می کند تا سازگاری آن در سیستم عامل های مختلف (ویندوز/لینوکس) تضمین گردد.

۲.۱.۷ نام کالکشن پایگاه داده Weaviate

```
WEAVIATE_COLLECTION_NAME = "Multimodal_Collection"
```

این نام برای کالکشنی است که داده‌های چندحالتی (شامل متن، تصویر، صدا) در پایگاه داده Weaviate در آن ذخیره می‌شوند و متناظر با همان *collection* ایجاد شده در مراحل پیشین است.

۳.۱.۷ بارگذاری متغیرهای محیطی از *env*.

```
load_dotenv()
```

این تابع به دنبال فایلی به نام *env*. در پوشه‌ی جاری یا والد آن جستجو کرده و تمام متغیرهای محیطی تعریف شده در آن را وارد محیط اجرایی می‌کند. فایل مذکور در مراحل اولیه پروژه ایجاد و کلید *API* لازم در آن تنظیم شده است.

۴.۱.۷ تنظیم دستگاه اجرا (CPU یا GPU)

```
DEVICE = "cuda" if torch.cuda.is_available() else "cpu"
```

- اگر سیستم مجهز به واحد پردازش گرافیکی (GPU) با پشتیبانی از CUDA باشد، از GPU استفاده خواهد شد.
- در غیر این صورت، عملیات بر روی واحد پردازش مرکزی (CPU) انجام می‌پذیرد.

۵.۱.۷ پیکربندی مدل CLIP

```
CLIP_MODEL_NAME = "ViT-B-32"
PRETRAINED_LOCAL_PATH =
    f"{PROJECT_ROOT_PATH}/open_clip_weights/ViT-B-32-openai/open_clip_model.safetensors"
```

- *CLIP_MODEL_NAME*: مدل *CLIP* منتخب که در اینجا *Vision Transformer* با اندازه‌ی *Base* و اندازه‌ی پچ ۳۲ (*ViT-B-32*) است.
- *PRETRAINED_LOCAL_PATH*: مسیر فایل وزن‌های از پیش‌آموزش‌دیده مدل است که به صورت محلی در پوشه‌ی پروژه ذخیره شده است. فایل *open_clip_model.safetensors* حاوی وزن‌های از پیش‌آموزش‌دیده مدل است.
- مدل *CLIP* به منظور تبدیل تصویر و متن به یک فضای برداری مشترک استفاده می‌گردد.

۶.۱.۷ پیکربندی مدل Whisper

```
WHISPER_MODEL = "small"
```

- این مدل از *OpenAI* برای تبدیل گفتار به متن (*Speech-to-Text*) به کار می‌رود.
- مقدار *WHISPER_MODEL* می‌تواند یکی از مقادیر *small*، *medium* یا *large* باشد.
- مدل‌های *medium* و *large* برای بارگذاری روی *CPU* مناسب نبوده و به *GPU* نیاز دارند.
- مدل *small* برای بارگذاری روی *CPU* مناسب است اما دقت کمتری نسبت به دو مدل دیگر دارد. این مدل برای داده‌های صوتی به زبان انگلیسی عملکرد مطلوبی دارد، اما برای زبان‌هایی غیر از انگلیسی (مانند فارسی) دارای عملکرد ضعیف‌تری است.

۷.۱.۷ پیکربندی مدل زبانی بزرگ (LLM)

```
LLM_MODEL_NAME = "openai/gpt-5-image-mini"
LLM_API_BASE = "https://openrouter.ai/api/v1"
LLM_API_KEY = os.getenv("LLM_API_KEY")
```

- *LLM_MODEL_NAME*: مدل زبانی مورد استفاده است که در اینجا نسخه‌ی کوچک از *GPT-5* با قابلیت پشتیبانی از تصویر می‌باشد.
- *LLM_API_BASE*: آدرس سرور *OpenRouter* برای فراخوانی *API* مدل‌ها.
- *LLM_API_KEY*: از فایل *env*. خوانده می‌شود و برای احراز هویت *API* الزامی است.

۲.۷ فایل `ai_models.py`

در این پروژه، مدل `CLIP` برای تولید بردارهای جاسازی (*embedding*) و مدل `Whisper` برای تبدیل صوت به متن استفاده شده است.

این دو مدل با استفاده از کد زیر بارگذاری شده‌اند. این کد در فایل `ai_models.py` واقع در پوشه `server` قرار دارد.

```
1 import whisper
2 import open_clip
3 from .config import WHISPER_MODEL, CLIP_MODEL_NAME, PRETRAINED_LOCAL_PATH, DEVICE
4
5 """ Whisper """
6 whisper_model = whisper.load_model(WHISPER_MODEL, device=DEVICE)
7
8 """ Open_CLIP """
9 clip_model, _, preprocess = open_clip.create_model_and_transforms(
10     CLIP_MODEL_NAME, pretrained=PRETRAINED_LOCAL_PATH
11 )
12 tokenizer = open_clip.get_tokenizer(CLIP_MODEL_NAME)
13 clip_model.to(DEVICE)
14 clip_model.eval()
```

این قطعه کد وظیفه‌ی بارگذاری و آماده‌سازی مدل‌های چندحالتی (صوت و تصویر) را برعهده دارد که در یک سیستم مبتنی بر بازیابی پیشرفته‌ی چندحالتی (مانند *Multimodal RAG Chatbot*) مورد استفاده قرار می‌گیرد. ساختار این فایل به گونه‌ای طراحی شده که مدل‌ها تنها یک‌بار بارگذاری شوند و در طول چرخه‌ی عمر پروژه آماده‌ی استفاده باشند.

۱.۲.۷ بارگذاری مدل `Whisper` (تبدیل گفتار به متن)

```
1 """ Whisper """
2 whisper_model = whisper.load_model(WHISPER_MODEL, device=DEVICE)
```

- `whisper.load_model()`: مدل انتخاب‌شده (`WHISPER_MODEL`) را بارگذاری می‌کند (مانند `small`, `base`, `medium`, `large`).
- پارامتر `device=DEVICE` تضمین می‌کند که مدل بر روی واحد پردازش گرافیکی (`GPU`) یا مرکزی (`CPU`) مناسب تنظیم شود.
- خروجی `whisper_model` یک شیء مدل است که قابلیت تبدیل ورودی صوتی به متن را فراهم می‌آورد.

۲.۲.۷ بارگذاری مدل `CLIP` (درک تصویر و متن)

```
1 """ Open_CLIP """
2 clip_model, _, preprocess = open_clip.create_model_and_transforms(
3     CLIP_MODEL_NAME, pretrained=PRETRAINED_LOCAL_PATH
4 )
```

- `open_clip.create_model_and_transforms()` سه مقدار را برمی‌گرداند:
 ۱. مدل `CLIP`: قابلیت تبدیل متن و تصویر به بردارهای مشترک در یک فضای برداری یکسان را دارد.
 ۲. مدیریت کلاس‌ها (`head`): که در اینجا با استفاده از `_` نادیده گرفته شده است زیرا مورد استفاده‌ی مستقیم قرار نمی‌گیرد.
 ۳. تابع `preprocess`: برای پیش‌پردازش تصویر (`normalization, resize` و غیره) پیش از ورود به مدل استفاده می‌شود.
- پارامتر `pretrained=PRETRAINED_LOCAL_PATH` باعث می‌شود وزن‌های مدل از فایل محلی (`safetensors`) بارگذاری شوند.

۳.۲.۷ بارگذاری مدل `Tokenizer` مدل `CLIP`

```
1 tokenizer = open_clip.get_tokenizer(CLIP_MODEL_NAME)
```

- توکنایزر (`Tokenizer`) مسئول تبدیل متن به ورودی عددی (`tokens`) قابل فهم برای مدل `CLIP` است.
- به عنوان مثال، برای توکن‌سازی یک متن: `text = tokenizer(["a dog running in the park"])`.
- در مدل `CLIP`، متن و تصویر هر دو به یک فضای برداری مشترک نگاشت می‌شوند تا امکان اندازه‌گیری میزان شباهت معنایی بین آن‌ها فراهم گردد.

```
1 clip_model.to(DEVICE)
2 clip_model.eval()
```

- مدل را مطابق با پیکربندی موجود در فایل *config* به GPU یا CPU منتقل می‌کند.
- *eval()*: مدل را در حالت استنتاج (*inference*) قرار می‌دهد، که منجر به:
 - غیرفعال شدن مکانیسم‌هایی نظیر *Dropout* و *BatchNorm* می‌شود.
 - آماده‌سازی مدل صرفاً برای پیش‌بینی و نه برای فرآیند آموزش.

۳.۷ فایل *database.py*

برای بازیابی (*retrieve*) فایل‌هایی که در مراحل پیشین به پایگاه داده (*database*) وارد شده‌اند، لازم است که به نمونه‌ی ایجاد شده از پایگاه داده متصل شویم. کد زیر، که در فایل *database.py* در پوشه‌ی *server* قرار دارد، همین وظیفه را بر عهده دارد:

```
1 try:
2     weaviate_client = weaviate.connect_to_local()
3     print(" Connected to Weaviate")
4 except Exception as e:
5     print(f" Weaviate connection failed: {e}")
6     weaviate_client = None
```

۱.۳.۷ توضیح دقیق عملکرد

اتصال به Weaviate محلی:

```
1 weaviate_client = weaviate.connect_to_local()
```

- این تابع تلاش می‌کند به نسخه‌ی محلی (*local*) از پایگاه داده‌ی وکتوری Weaviate که بر روی سیستم یا از طریق *Docker* در حال اجرا است، متصل شود.
- به‌صورت پیش‌فرض، آدرس اتصال *http://localhost:8080* در نظر گرفته شده است.
- در صورت موفقیت‌آمیز بودن اتصال، یک شیء کلاینت (*client object*) برگردانده می‌شود که امکان تعامل با پایگاه داده را فراهم می‌سازد.

۴.۷ فایل *main.py*

پیش از اجرای پروژه، لازم است اطمینان حاصل شود که داکر (*Docker*) بر روی سیستم شما در حال اجرا است، زیرا اتصال به پایگاه داده (*database*) مستلزم فعال بودن داکر است. برای اجرای پروژه، یک واسط خط فرمان (*cmd*) باز کنید، به پوشه‌ی ریشه‌ی پروژه بروید و دستور زیر را اجرا نمایید:

```
1 python -m server.main
```

سپس، یک مرورگر وب باز کرده و به آدرس زیر مراجعه نمایید:

```
1 localhost:8000
```

با اجرای دستور *python -m server.main*، فایل اصلی پروژه یعنی *main.py* که در پوشه‌ی *server* قرار دارد، اجرا می‌شود.

۱.۴.۷ کد فایل *main.py*

کد زیر، محتوای کامل فایل *main.py* را نشان می‌دهد:

```
1 import os
2 from fastapi import FastAPI
3 from fastapi.middleware.cors import CORSMiddleware
4 from fastapi.responses import HTMLResponse
5 from fastapi.staticfiles import StaticFiles
6 from .search_routes import router as search_router
7
```

```

8 app = FastAPI(title="Multimodal RAG Chatbot")
9
10 # CORS
11 app.add_middleware(
12     CORSMiddleware,
13     allow_origins=["*"],
14     allow_credentials=True,
15     allow_methods=["*"],
16     allow_headers=["*"],
17 )
18
19 # Static mounts
20 app.mount("/client", StaticFiles(directory="client"), name="client")
21 app.mount("/content", StaticFiles(directory="content"), name="content")
22
23 # Routes
24 app.include_router(search_router)
25
26 @app.get("/", response_class=HTMLResponse)
27 async def home():
28     # Serve the index.html file properly
29     html_path = os.path.join("client", "index.html")
30     if os.path.exists(html_path):
31         with open(html_path, encoding="utf-8") as f:
32             return HTMLResponse(f.read())
33     else:
34         return HTMLResponse("<h1>Frontend not found</h1>", status_code=404)
35
36 if __name__ == "__main__":
37     import uvicorn
38     uvicorn.run(app, host="0.0.0.0", port=8000)

```

۲.۴.۷ توضیح کد

ساخت اپلیکیشن *FastAPI*:

```

1 app = FastAPI(title="Multimodal RAG Chatbot")

```

یک شیء *FastAPI* با عنوان مشخص شده ساخته می‌شود. این شیء نمایانگر کل برنامه‌ی وب (بک اند) پروژه است. تنظیم *CORS*:

```

1 app.add_middleware(
2     CORSMiddleware,
3     allow_origins=["*"],
4     allow_credentials=True,
5     allow_methods=["*"],
6     allow_headers=["*"],
7 )

```

این بخش مجوزهای لازم را برای دسترسی همه‌ی درخواست‌ها (*from any domain*) به *API* فراهم می‌کند. جزئیات این تنظیمات به شرح زیر است:

- *allow_origins=["*"]*: تمامی دامنه‌ها مجاز به ارسال درخواست هستند.
- *allow_credentials=True*: ارسال کوکی‌ها در درخواست‌ها مجاز است.
- *allow_methods=["*"]*: تمامی متدهای *HTTP* (شامل *GET*، *POST*، و غیره) مجازند.
- *allow_headers=["*"]*: تمامی سرآیندهای (*headers*) ارسالی مجاز هستند.

سرو فایل‌های استاتیک

```

1 app.mount("/client", StaticFiles(directory="client"), name="client")
2 app.mount("/content", StaticFiles(directory="content"), name="content")

```

با استفاده از این دو خط، دو مسیر مجزا برای فایل‌های استاتیک تعریف شده است:

- محتوای پوشه‌ی *client* در مسیر */client* قابل دسترسی است.
- محتوای پوشه‌ی *content* در مسیر */content* قابل دسترسی است.

افزودن مسیرهای دیگر از فایل جداگانه

```
1 app.include_router(search_router)
```

این دستور کلیه‌ی مسیرهای تعریف‌شده در فایل `search_routes.py` را به ساختار اصلی اپلیکیشن اضافه می‌کند.
تعریف مسیر اصلی (/)

```
1 @app.get("/", response_class=HTMLResponse)
2 async def home():
3     html_path = os.path.join("client", "index.html")
4     if os.path.exists(html_path):
5         with open(html_path, encoding="utf-8") as f:
6             return HTMLResponse(f.read())
7     else:
8         return HTMLResponse("<h1>Frontend not found</h1>", status_code=404)
```

در هنگام مراجعه‌ی کاربر به مسیر ریشه (/)، عملکرد به صورت زیر است:

- برنامه به دنبال فایل `client/index.html` می‌گردد.
- اگر فایل مورد نظر یافت شود، محتوای آن به عنوان پاسخ `HTML` بازگردانده و نمایش داده می‌شود.
- در غیر این صورت، پیام خطای `Frontend not found` به همراه کد وضعیت `404` بازگردانده خواهد شد.

اجرای برنامه با `Uvicorn`

```
1 if __name__ == "__main__":
2     import uvicorn
3     uvicorn.run(app, host="0.0.0.0", port=8000)
```

این بلوک تضمین می‌کند که در صورت اجرای مستقیم فایل (`python main.py`):

- سرور `FastAPI` توسط `Uvicorn` اجرا گردد.
- سرور بر روی تمامی `IP`های محلی (`0.0.0.0`) و پورت `8000` در دسترس باشد.

۳.۴.۷ قابلیت‌های جستجو و درخواست‌ها

پس از دسترسی به آدرس `localhost:8000` در مرورگر، می‌توان درخواست‌هایی با ترکیب‌های چندحالتی (`multimodal`) زیر را به سیستم ارسال کرد. درخواست‌ها می‌توانند شامل هر ترکیبی از متن، صوت و تصویر باشند:

۱. فقط متن
۲. فقط تصویر
۳. فقط فایل صوتی
۴. متن + تصویر
۵. متن + فایل صوتی
۶. تصویر + فایل صوتی
۷. متن + تصویر + فایل صوتی

توجه: ورود تعداد زیادی فایل صوتی یا تصویری به طور همزمان توصیه نمی‌شود، هرچند محدودیتی در تعداد وجود ندارد. دلایل عدم توصیه به شرح زیر است:

- این اقدام می‌تواند منجر به کاهش کارایی پاسخ‌دهی سیستم (`response`) گردد (که جزئیات آن در ادامه بررسی خواهد شد).
- تمام داده‌های تصویری به مدل زبان بزرگ (`LLM`) ارسال می‌شوند. در صورت زیاد بودن حجم داده‌ها، احتمال اتمام توکن‌های مجاز برای استفاده‌ی رایگان از مدل وجود دارد و `LLM` برای مدتی پاسخگو نخواهد بود.

۵.۷ فایل `utils.py` و `search_routes.py`

هنگامی که یک درخواست جستجو در بخش فرانت (*Frontend*) ایجاد می‌شود، از طریق جاوا اسکریپت (*JavaScript*) به مسیر *API* زیر درخواست زده می‌شود:

```
1 @router.post("/multimodal")
2 async def search_multimodal(
3     query: str = Form(default=""), files: List[UploadFile] = File(default=[])
4 ):
5     pass
```

این تابع در فایل `search_routes.py` واقع در پوشه‌ی `server` تعریف شده است.

۱.۵.۷ توضیح کلی عملکرد مسیر `/multimodal`

این مسیر یک درخواست چندمودال (*multimodal*) را دریافت می‌کند (شامل متن و/یا فایل‌های آپلودشده مانند تصویر/صوت). ابتدا ورودی‌ها پردازش شده و به بردارهای جاسازی (*embedding*) تبدیل می‌شوند. سپس یک بردار نهایی (تک یا میانگین بردارها) ساخته شده، با آن در پایگاه داده‌ی وکتوری (*Vector Database*) جستجو انجام می‌شود. در نهایت، نتایج نرمال‌سازی شده، داده‌های لازم برای مدل زبان بزرگ (*LLM*) آماده و به آن فرستاده می‌شود و پاسخ *JSON* شامل نتایج بازیابی شده و خروجی *LLM* به کلاینت بازگردانده می‌شود. در انتهای فرآیند، فایل‌های موقتی حذف می‌گردند.

۲.۵.۷ بررسی گام به گام کد

۱. امضای تابع و ورودی‌ها

```
1 @router.post("/multimodal")
2 async def search_multimodal(
3     query: str = Form(default=""),
4     files: List[UploadFile] = File(default=[])
5 ):
6     pass
```

- این یک مسیر *POST* است که یک فیلد فرم *query* (متن آزاد) و لیستی از فایل‌های آپلودشده (*files*) را دریافت می‌کند.
- *UploadFile* یک کلاس استاندارد از *FastAPI* است که برای مدیریت فایل‌های آپلودی استفاده می‌شود.
- تابع به صورت ناهمگام (*async*) تعریف شده است تا امکان استفاده از عملیات‌های غیرهمزمان (مانند `await save_temp_file(file)`) فراهم گردد.

۲. آماده‌سازی متغیرها

- *temp_paths*: مسیر فایل‌های موقتی که برای پاکسازی نهایی روی دیسک ذخیره شده‌اند.
- *embeddings*: لیستی از بردارهای *embedding* استخراج‌شده از متن، تصویر یا صوت.
- *response_data*: دیکشنری خروجی که در نهایت به فرانت بازگردانده می‌شود.

۳. تشخیص وجود متن و پردازش آن اگر درخواست ورودی حاوی متن باشد، بخش زیر اجرا می‌شود:

```
1 has_text = bool(query and query.strip())
2
3 if has_text:
4     english_query = normalize_text_to_english(query)
5     text_emb = embed_text(english_query)
6     embeddings.append(text_emb)
7     response_data["original_query"] = query
8     response_data["processed_query"] = english_query
```

- از تابع `normalize_text_to_english` برای ترجمه‌ی متن استفاده می‌شود. در نتیجه اگر متن ورودی به زبانی غیر انگلیسی باشد، ابتدا به انگلیسی ترجمه می‌گردد.
- سپس برای متن انگلیسی‌شده، با استفاده از متد `embed_text` یک بردار جاسازی (*embedding*) به دست می‌آید.
- *embedding* متن به لیست *embeddings* اضافه می‌شود و اطلاعات اصلی و پردازش‌شده در *response_data* ذخیره می‌گردد.

۴. پردازش فایل‌های آپلودی (ذخیره موقت و استخراج embedding)

```

1 for file in files:
2     temp_path = await _save_temp_file(file)
3     temp_paths.append(temp_path)
4
5     if file.content_type.startswith("image/"):
6         img_emb = embed_image(temp_path)
7         embeddings.append(img_emb)
8         try:
9             img_obj = Image.open(temp_path).copy()
10        except Exception as e:
11            img_obj = None
12            image_files.append(img_obj)
13
14    elif file.content_type.startswith("audio/"):
15        audio_emb, audio_meta = _process_audio(temp_path)
16        embeddings.append(audio_emb)
17        audio_transcriptions.append(audio_meta)

```

- هر فایل ابتدا با `save_temp_file` به صورت موقت ذخیره شده و مسیر آن برای پاکسازی نهایی نگهداری می‌شود.
- **برای تصاویر:** تابع `embed_image` بردار تصویر را ساخته و به لیست اضافه می‌کند. همچنین، تصویر با کتابخانه‌ی `PIL` باز شده و در `image_files` نگهداری می‌شود تا برای ارسال به `LLM` یا پیش‌نمایش در دسترس باشد.
- **برای صدا:** ابتدا با استفاده از تابع `process_audio`، بردار جاسازی و اطلاعات فراداده‌ای (`metadata`) استخراج و به لیست‌های مربوطه اضافه می‌شوند. این فراداده حاوی نام فایل، متن فایل به زبان اصلی، متن فایل به زبان انگلیسی و زبان شناسایی‌شده‌ی فایل است که در ادامه توضیح داده خواهد شد.

۵. به دست آوردن متن فایل‌های صوتی

- اگر تنها یک فایل صوتی وجود داشته و هیچ متن ورودی‌ای توسط کاربر ارائه نشده باشد، مستقیماً نتیجه‌ی `transcript` (`transcription`) در `response_data` قرار می‌گیرد.
- در غیر این صورت، تمامی `transcription` ها در قالب یک لیست بازگردانده می‌شوند.

۶. تعیین بردار نهایی جستجو (single vs multimodal mean)

```

1 if not embeddings:
2     return JsonResponse({"error": "No valid text or media found to process."},
3                           status_code=400)
4
5 if len(embeddings) == 1:
6     q_emb = embeddings[0]
7     response_data["search_mode"] = "single"
8 else:
9     q_emb_array = np.array(embeddings)
10    mean_emb = np.mean(q_emb_array, axis=0)
11    q_emb = mean_emb
12    response_data["search_mode"] = f"multimodal_mean ({len(embeddings)} inputs)"

```

- اگر هیچ بردار جاسازی (`embedding`) تولید نشده باشد، خطای 400 بازگردانده می‌شود.
- اگر تنها یک بردار وجود داشته باشد، همان به عنوان بردار جستجوی نهایی (`q_emb`) تعیین می‌شود (`single mode`).
- اگر چندین بردار وجود داشته باشد، میانگین این بردارها گرفته می‌شود تا یک بردار واحد چندمودال ایجاد گردد.
- نکته: این رویکرد (میانگین‌گیری) ساده‌ترین روش ترکیب مودال‌ها است. مشکل آن این است که اگر ورودی‌ها دارای مفاهیم دور از هم باشند، بردار میانگین ممکن است دقت بازیابی مناسبی نداشته باشد. این بردار نهایی (`q_emb`) برای جستجو در پایگاه داده استفاده خواهد شد.

۷. اجرای جستجو در پایگاه داده برداری با استفاده از بردار نهایی `q_emb`، جستجو برای هر مودال به صورت جداگانه اجرا می‌شود:

```

1 text_res = search_by_embedding(q_emb, "text", top_k=5)
2 image_res = search_by_embedding(q_emb, "image", top_k=5)
3 audio_res = search_by_embedding(q_emb, "audio", top_k=5)

```

- تابع `search_by_embedding` فراخوانی می‌شود. این تابع بردار نهایی را دریافت کرده و بر اساس مدل‌بسته مشخص شده `(audio, image, text)`، تعداد 5 داده‌ی مشابه را از پایگاه داده بازیابی می‌کند.
- این فرآیند تضمین می‌کند که از هر نوع مدل‌بسته، 5 داده‌ی نزدیک به ورودی کاربر بازیابی گردد.

۸. نرمال‌سازی نتایج و آماده‌سازی برای LLM سپس با استفاده از تابع `normalize_results`، نتایج به فرمت یکپارچه برای ارسال به LLM و فرانت تبدیل و در `response_data` ذخیره می‌شوند.

```
1 normalized_texts = _normalize_results(text_res)
2 normalized_images = _normalize_results(image_res)
3 normalized_audios = _normalize_results(audio_res)
4
5 response_data.update({
6     "text_results": normalized_texts,
7     "image_results": normalized_images,
8     "audio_results": normalized_audios,
9 })
```

۹. فیلتر کردن و آماده‌سازی داده‌های نهایی برای LLM در این مرحله، ما دارای موارد زیر هستیم:

۱. تمام کوئری‌های کاربر (متن، تصویر، صوت): کوئری متن در متغیر `query`، کوئری‌های تصویر در لیست `image_files` و متن تبدیل‌شده صوت در لیست `audio_transcriptions` ذخیره شده‌اند.
۲. یک لیست 5 تایی از متون بازیابی‌شده (`normalized_texts`) که از دیتابیس بازیابی و نرمال شده‌اند.
۳. یک لیست 5 تایی از تصاویر بازیابی‌شده (`normalized_images`) که از دیتابیس بازیابی و نرمال شده‌اند.
۴. یک لیست 5 تایی از متون مربوط به صوت‌های بازیابی‌شده (`normalized_audios`) که از دیتابیس بازیابی و نرمال شده‌اند.

برای جلوگیری از ارسال حجم زیاد داده به LLM و مدیریت توکن‌ها، تنها زیرمجموعه‌ای از داده‌ها (کوئری‌های کاربر و فقط یک مورد از هر نوع داده بازیابی شده) در دیکشنری `llm_data` ذخیره و به LLM ارسال می‌شوند. داده‌هایی که در دیکشنری `llm_data` ذخیره شده و به LLM فرستاده می‌شوند عبارت‌اند از:

- تمام داده‌های کاربر.
 - یک مورد از متون بازیابی‌شده (`normalized_texts[0]`).
 - یک مورد از تصاویر بازیابی‌شده (`normalized_images[0]`).
 - یک مورد از متون صوتی بازیابی‌شده (`normalized_audios[0]`).
- همچنین، یک نسخه دیگر از این داده‌ها به نام `llm_data_modified_for_front` ساخته می‌شود که همان محتویات داده‌های بازیابی‌شده را دارد اما **بدون کوئری‌های کاربر**، که صرفاً برای راحتی کار در فرانت (`Frontend`) ایجاد شده است.

```
1 llm_data = {
2     "user_text_query": query if has_text else None,
3     "user_image_queries": image_files if image_files else None,
4     "user_audio_queries": (
5         [a.get("processed_text") for a in audio_transcriptions]
6         if audio_transcriptions
7         else None
8     ),
9     "retrieved_texts": (
10        [t.get("content") for t in normalized_texts[:1]]
11        if normalized_texts
12        else None
13    ),
14    "retrieved_images": (
15        [
16            (
17                Image.open(PROJECT_ROOT_PATH + i.get("filePath"))
18                if i.get("filePath")
19                else (
20                    Image.open(requests.get(i.get("url"), stream=True).raw)
21                    if i.get("url")
22                    else None
23                )
24            )
```



```

25         for i in normalized_images[:1]
26     ]
27     if normalized_images
28     else None
29 ),
30 "retrieved_audios": (
31     [
32         a.get("content") or a.get("processed_text")
33         for a in normalized_audios[:1]
34     ]
35     if normalized_audios
36     else None
37 ),
38 }
39
40 llm_data_modified_for_front = {
41     "retrieved_texts": (
42         [t.get("content") for t in normalized_texts[:1]]
43         if normalized_texts
44         else None
45     ),
46     "retrieved_images": (
47         [i.get("filePath") or i.get("url") for i in normalized_images[:1]]
48         if normalized_images
49         else None
50     ),
51     "retrieved_audios": (
52         [a.get("filePath") for a in normalized_audios[:1]]
53         if normalized_audios
54         else None
55     ),
56 }

```

۱۰. ارسال داده به LLM و بازگرداندن پاسخ دیکشنری `llm_data` به تابع `feed_data_into_llm` ارسال می‌شود. خروجی این تابع که پاسخ متنی LLM است، در `response_data["llm_response"]` برای نمایش در فرانت ذخیره می‌شود.

```

1 response_data["llm_response"] = feed_data_into_llm(llm_data)

```

۳.۵.۷ توابع کمکی (Utility Functions) در `utils.py`

در فرآیند بالا، از توابع کمکی متعددی استفاده شده است که در ادامه توضیح داده می‌شوند:

۱. تابع `normalize_text_to_english` هدف: نرمال‌سازی متن ورودی به زبان انگلیسی از طریق ترجمه.

```

1 def normalize_text_to_english(text: str) -> str:
2     """
3     Convert any input text into English if it is not already English.
4     Works offline for language detection + uses Google Translate API.
5     """
6     if not text or text.strip() == "":
7         return text
8
9     try:
10         translated = GoogleTranslator(source='auto', target='en').translate(text)
11         return translated
12     except Exception as e:
13         print(f"Translation error: {e}")
14         return text

```

- این تابع از `GoogleTranslator` (از پکیج `deep-translator`) استفاده می‌کند.

- `source='auto'` زبان ورودی را به صورت خودکار تشخیص می‌دهد.

- `target='en'` تضمین می‌کند که خروجی همیشه به زبان انگلیسی باشد.

۲. تابع `embed_text` هدف: تبدیل یک رشته متنی به یک بردار عددی (`embedding`) با استفاده از مدل `CLIP`.

```

1 def _to_numpy(t: torch.Tensor) -> np.ndarray:
2     return t.detach().cpu().numpy().reshape(-1)
3
4 def embed_text(text: str) -> np.ndarray:

```

```

5 tokens = tokenizer([text]).to(DEVICE)
6 with torch.no_grad():
7     text_features = clip_model.encode_text(tokens)
8     text_features /= text_features.norm(dim=-1, keepdim=True)
9     return _to_numpy(text_features)

```

- `tokens = tokenizer([text]).to(DEVICE)`: متن را به `tokens` تبدیل کرده و به دستگاه مناسب (CPU/GPU) منتقل می‌کند.
 - `with torch.no_grad()`: برای صرفه‌جویی در منابع و سرعت، محاسبه‌ی گرادیان‌ها غیرفعال می‌شود.
 - `text_features = clip_model.encode_text(tokens)`: ویژگی‌های متنی توسط مدل `CLIP` استخراج می‌شوند.
 - نرمال‌سازی بردارها: طول (`norm`) هر بردار به 1 نرمال می‌شود تا برای مقایسه‌ی شباهت کسینوسی (`cosine similarity`) مناسب باشد.
 - `return _to_numpy(text_features)`: تانسور `PyTorch` به آرایه‌ی `NumPy` تبدیل می‌گردد.
۳. تابع `embed_image` هدف: تبدیل یک تصویر به یک بردار عددی (`embedding`) با استفاده از مدل `CLIP`.

```

1 def embed_image(image_path: str) -> np.ndarray:
2     img = Image.open(image_path).convert("RGB")
3     x = preprocess(img).unsqueeze(0).to(DEVICE)
4     with torch.no_grad():
5         image_features = clip_model.encode_image(x)
6         image_features /= image_features.norm(dim=-1, keepdim=True)
7     return _to_numpy(image_features)

```

- `img = Image.open(image_path).convert("RGB")`: تصویر با `Pillow` باز و به فرمت `RGB` تبدیل می‌شود.
 - `x = preprocess(img).unsqueeze(0).to(DEVICE)`: تصویر توسط تابع `preprocess` (دریافت شده از `open_clip`) پیش‌پردازش می‌شود و به دستگاه منتقل می‌گردد. `unsqueeze(0)` بعد `batch` را اضافه می‌کند.
 - `image_features = clip_model.encode_image(x)`: ویژگی‌های تصویری توسط مدل `CLIP` استخراج می‌شوند.
 - بردار ویژگی‌ها نرمال‌سازی شده و در نهایت به آرایه‌ی `NumPy` تبدیل می‌شود.
۴. تابع `save_temp_file` هدف: ذخیره‌ی فایل آپلودشده (`UploadFile`) در یک مسیر موقت و بازگرداندن مسیر آن.

```

1 async def save_temp_file(file: UploadFile) -> str:
2     """Saves UploadFile to a temp path and returns the path."""
3     # Ensure a file extension, default to .tmp if none
4     suffix = os.path.splitext(file.filename)[1] or ".tmp"
5     with tempfile.NamedTemporaryFile(delete=False, suffix=suffix) as tf:
6         tf.write(await file.read())
7     return tf.name

```

- تابع `asynchronous` است و از `await` برای خواندن محتوای فایل استفاده می‌کند.
 - با `tempfile.NamedTemporaryFile(delete=False, ...)`، فایل موقت با پسوند مناسب ایجاد می‌شود و تضمین می‌گردد که پس از بسته شدن نیز (`delete=False`) روی دیسک باقی بماند.
۵. تابع `process_audio` هدف: تبدیل صوت به متن، شناسایی زبان، ترجمه‌ی آن به انگلیسی (در صورت لزوم)، تولید `embedding` و جمع‌آوری فراداده.

```

1 def process_audio(temp_path: str) -> Tuple[List[float], dict]:
2     """
3     Transcribes audio, normalizes text, embeds it,
4     and returns (embedding, metadata_dict).
5     """
6     # 1. Whisper -> Text + Language Detection
7     whisper_res = whisper_model.transcribe(temp_path)
8     detected_text = whisper_res.get("text", "").strip()
9     detected_lang = whisper_res.get("language", "")
10
11     # 2. If Persian -> Translate to English
12     if detected_lang.startswith("fa"):

```

```

13 processed_text = normalize_text_to_english(detected_text)
14 else:
15     processed_text = detected_text
16
17 # 3. Embed transcribed text
18 q_emb = embed_text(processed_text)
19
20 metadata = {
21     "file_name": os.path.basename(temp_path),
22     "detected_text": detected_text,
23     "processed_text": processed_text,
24     "detected_language": detected_lang,
25 }
26
27 return q_emb, metadata

```

تعریف تابع ورودی/خروجی

```

1 def process_audio(temp_path: str) -> Tuple[List[float], dict]:

```

- ورودی: *temp_path: str* - مسیر فایل صوتی (مثل *wav* یا *mp3*) که قبلاً روی دیسک ذخیره شده است.
- خروجی: یک تاپل (*tuple*) شامل دو مقدار:

- *q_emb: List[float]* - بردار جاسازی (*embedding*) متن آن فایل صوتی.
 - *metadata: dict* - دیکشنری شامل اطلاعات اضافی درباره‌ی فایل و متن تشخیص داده‌شده.

گام به گام اجرای تابع

۱. تبدیل صوت به متن با *Whisper* و شناسایی زبان

```

1 whisper_res = whisper_model.transcribe(temp_path)
2 detected_text = whisper_res.get("text", "").strip()
3 detected_lang = whisper_res.get("language", "")
4

```

- *whisper_model.transcribe(temp_path)* فایل صوتی را به مدل (*Whisper*) (*OpenAI Speech Recognition*) ارسال می‌کند تا متن و زبان را استخراج کند.
- *detected_text*: متن استخراج شده از صوت.
- *detected_lang*: کد زبان شناسایی شده (مانند "*en*" یا "*fa*").

نتیجه این بخش: صوت تبدیل به متن و شناسایی زبان انجام می‌شود.

۲. ترجمه به انگلیسی (در صورت فارسی بودن زبان)

```

1 if detected_lang.startswith("fa"):
2     processed_text = normalize_text_to_english(detected_text)
3 else:
4     processed_text = detected_text
5

```

- اگر زبان تشخیص داده شده با "*fa*" شروع شود (فارسی)، تابع *normalize_text_to_english()* فراخوانی می‌شود تا متن به انگلیسی ترجمه گردد.
- این کار متن نهایی را برای تولید *embedding* سازگار با مدل چندمدال (*CLIP*) آماده می‌کند.

نتیجه این بخش: متن نهایی (*processed_text*) به زبان انگلیسی آماده می‌شود.

۳. گرفتن *embedding* از متن

```

1 q_emb = embed_text(processed_text)
2

```

- تابع *embed_text()* (که قبلاً توضیح داده شد) متن پردازش شده را به یک بردار عددی (*embedding*) تبدیل می‌کند که نشان‌دهنده معنای جمله است.

نتیجه این بخش: متن نهایی به فضای برداری (*Vector Space*) نگاشته می‌شود.

۴. ساخت *metadata*

```
1 metadata = {  
2     "file_name": os.path.basename(temp_path),  
3     "detected_text": detected_text,  
4     "processed_text": processed_text,  
5     "detected_language": detected_lang,  
6 }  
7
```

توضیحات فیلدهای *metadata* دیکشنری متادیتا شامل اطلاعات زیر برای ثبت و استفاده‌ی بعدی است:

- "file_name": نام فایل صوتی.
- "detected_text": متن خام تشخیص داده‌شده توسط مدل *Whisper*.
- "processed_text": متن نهایی بعد از ترجمه به انگلیسی یا نرمال‌سازی (متنی که برای تولید *embedding* استفاده شده است).
- "detected_language": کد زبان شناسایی‌شده.

نتیجه این بخش: همه اطلاعات متنی و زبانی همراه با نام فایل در قالب *metadata* جمع‌آوری می‌شوند.

۵. برگرداندن خروجی

```
1 return q_emb, metadata  
2
```

- خروجی شامل *Embedding (q_emb)* برای جستجو در دیتابیس برداری و *Metadata (metadata)* برای اطلاعات جانبی است.

۶. توابع *get_embedding* و *search_by_embedding* هدف: تولید *embedding* (نمایش برداری) برای کوئری ورودی و اجرای جستجوی تشابه برداری (*Vector Similarity Search*) در پایگاه داده *Weaviate*. این دو تابع در کنار یکدیگر کار می‌کنند تا یک کوئری ورودی (متن یا تصویر) را به بردار تبدیل کرده و نزدیک‌ترین موارد مشابه را از دیتابیس بازیابی کنند.

```
1 def get_embedding(modality: str, data: str):  
2     if modality == "text":  
3         return embed_text(data)  
4     elif modality == "image":  
5         return embed_image(data)  
6     raise ValueError("modality must be 'text' or 'image'")  
7  
8 def search_by_embedding(query_embedding, modality: str, top_k=3):  
9     try:  
10        collection = weaviate_client.collections.get(WEAVIATE_COLLECTION_NAME)  
11        results = collection.query.near_vector(  
12            near_vector=query_embedding.tolist(),  
13            limit=top_k,  
14            filters=weaviate.classes.query.Filter.by_property("modality").equal(  
15                modality  
16            ),  
17        )  
18        found = []  
19        if results.objects:  
20            for obj in results.objects:  
21                found.append({"properties": obj.properties})  
22        return found  
23    except Exception as e:  
24        print(f"Error searching {modality}: {e}")  
25        return []
```

تابع *get_embedding(modality: str, data: str)* وظیفه: بر اساس نوع داده (*modality*: "text" یا "image")، *embedding* مناسب را تولید و برمی‌گرداند.

- *if modality == "text": return embed_text(data)*: اگر نوع داده متنی باشد، متن به یک بردار عددی تبدیل می‌شود.

• `elif modality == "image": return embed_image(data)`: اگر نوع داده تصویری باشد، فایل تصویر به یک `embedding` عددی تبدیل می‌شود.

• `raise ValueError(...)` ...: در صورت نامعتبر بودن `modality`، یک خطا پرتاب می‌شود.

خروجی: یک آرایه‌ی عددی (`numpy array`) شامل `embedding` کوثری ورودی.
تابع `search_by_embedding(query_embedding, modality: str, top_k=3)` وظیفه: جستجوی شباهت برداری با استفاده از `embedding` کوثری در پایگاه داده `Weaviate` و بازگرداندن نزدیک‌ترین نتایج.

۱. اتصال به کالکشن (`Collection`):

```
1 collection = weaviate_client.collections.get(WEAVIATE_COLLECTION_NAME)
```

به کالکشن مورد نظر در `Weaviate` که حاوی داده‌های چندمودال است، متصل می‌شود.

۲. اجرای جستجوی برداری:

```
1 results = collection.query.near_vector(  
2     near_vector=query_embedding.tolist(),  
3     limit=top_k,  
4     filters=weaviate.classes.query.Filter.by_property("modality").equal(modality),  
5 )  
6
```

• `collection.query.near_vector`: متد اصلی برای جستجوی تشابه برداری.

• `limit=top_k`: تعداد نتایج مورد نیاز را مشخص می‌کند (به طور پیش فرض ۳).

• `filters`: این فیلتر ضروری تضمین می‌کند که جستجو فقط در میان داده‌هایی انجام شود که نوع مدالیته‌ی آنها با کوثری ورودی یکسان است (مثلاً `embedding` یک متن، فقط در میان `embedding` متون جستجو شود).

۳. استخراج نتایج:

```
1 found = []  
2 if results.objects:  
3     for obj in results.objects:  
4         found.append({"properties": obj.properties})  
5 return found  
6
```

• داده‌ها و فراداده‌های (`metadata`) مرتبط با نزدیک‌ترین بردارها، از فیلد `obj.properties` استخراج شده و در لیست `found` ذخیره می‌شوند.

۴. مدیریت خطا:

```
1 except Exception as e:  
2     print(f"Error searching {modality}: {e}")  
3     return []  
4
```

در صورت بروز خطا در حین جستجو در `Weaviate`، پیام خطا چاپ شده و یک لیست خالی برگردانده می‌شود.

۷. تابع `normalize_results` هدف: ساده‌سازی خروجی جستجوهای `Weaviate` به یک لیست ساده از ویژگی‌ها (`properties`).

```
1 def normalize_results(results_list):  
2     """Converts Weaviate result list to a simple list of properties."""  
3     out = []  
4     for r in results_list or []:  
5         props = r.get("properties") if isinstance(r, dict) else None  
6         if props:  
7             out.append(props)  
8     return out
```

این تابع برای تبدیل نتایج خام دیتابیس به فرمت یکنواخت برای استفاده‌ی آسان‌تر طراحی شده است.

۶.۷ فایل llm.py

در نهایت، جهت تولید پاسخ توسط مدل زبان بزرگ (LLM)، از تابع `feed_data_into_llm` که در اسکریپت `llm.py` تعریف شده، استفاده می‌گردد.

۱.۶.۷ راه‌اندازی کلاینت OpenAI

در ابتدای فایل `llm.py`، کد زیر جهت برقراری ارتباط با مدل‌های زبان بزرگ، تعبیه شده است:

```
1 # === Initialize OpenAI Client ===
2 client = OpenAI(
3     base_url=LLM_API_BASE,
4     api_key=LLM_API_KEY
5 )
```

این قسمت از کد وظیفه راه‌اندازی (Initialization) کلاینت `OpenAI` را بر عهده دارد تا امکان استفاده از مدل‌های زبانی (LLM) مانند `GPT` فراهم گردد. متغیر `client` به عنوان رابط اصلی برای تعامل با `LLM` عمل می‌کند.

۱. `client = OpenAI(...)`: در این خط، یک شیء (object) از کلاس `OpenAI` ایجاد شده و در متغیر `client` ذخیره می‌شود. این شیء در واقع کلاینتی است که از طریق آن می‌توان درخواست‌ها را به مدل‌های `OpenAI` یا سرورهای سازگار با `API` آن (مانند سرورهای محلی یا `third-party LLM servers`) ارسال نمود.

۲. پارامترهای کلیدی:

- `base_url=LLM_API_BASE`: این پارامتر آدرس `URL` سروری که درخواست‌ها به آن ارسال می‌شود را مشخص می‌کند.

- در استفاده از `API` رسمی `OpenAI`: `LLM_API_BASE = "https://api.openai.com/v1"`
- در صورت استفاده از مدل‌های محلی (`local`) یا میزبانی‌شده توسط کاربر (مانند `self-hosted`) مانند `Ollama`، `Studio`، یا `LLM`: `LLM_API_BASE=(vLLM می‌تواند آدرس API محلی باشد، مثلاً: "http://localhost:8000/v1"`

- `api_key=LLM_API_KEY`: این کلید امنیتی جهت احراز هویت (authentication) کاربر در سرور الزامی است. مقدار آن معمولاً یک رشته محرمانه مانند `"sk-xxxx..."` است.

۳. نتیجه: پس از اجرای این خطوط، متغیر `client` به یک شیء فعال از نوع `OpenAI` تبدیل می‌شود.

لازم به ذکر است که مقادیر `LLM_API_KEY` و `LLM_API_BASE` از فایل تنظیمات `config.py` بارگذاری می‌شوند.

۲.۶.۷ تشریح تابع feed_data_into_llm

```
1 def feed_data_into_llm(llm_data: dict) -> str:
2     """Send multimodal data (text + image + audio) to OpenRouter model."""
3
4     def section(title: str, content: str) -> str:
5         return f"\n\n### {title}\n\n{content.strip()}"
6
7     sections = []
8
9     # --- User Text ---
10    if llm_data.get("user_text_query"):
11        sections.append(section("User Text Query", llm_data["user_text_query"]))
12
13    # --- User Images ---
14    image_base64s = []
15    if llm_data.get("user_image_queries"):
16        for img in llm_data["user_image_queries"]:
17            pil_img = to_pil_image(img)
18            if pil_img:
19                pil_img = pil_img.resize((128, 128), Image.LANCZOS)
20                image_base64s.append(image_to_base64_data_url(pil_img))
21        sections.append(section("User Image Queries",
22                                f"{len(image_base64s)} image(s) attached." if image_base64s else
23                                "<no usable images>"))
24
25    # --- User Audio ---
26    if llm_data.get("user_audio_queries"):
27        aud_list = "\n\n".join(
28            [f"Audio {i+1} (Transcription):\n{t}" for i, t in
29             enumerate(llm_data["user_audio_queries"])]
30        )
```

```

30     )
31     sections.append(section("User Audio Queries", aud_list))
32
33     # --- Retrieved Texts ---
34     if llm_data.get("retrieved_texts"):
35         txt_list = "\n\n".join(
36             [f"Retrieved Text {i+1}:\n{t}" for i, t in
37              enumerate(llm_data["retrieved_texts"])]
38         )
39         sections.append(section("Retrieved Texts", txt_list))
40
41     # --- Retrieved Images ---
42     retrieved_image_base64s = []
43     print("Retrieved images count:", len(llm_data.get("retrieved_images", [])))
44     if llm_data.get("retrieved_images"):
45         for img in llm_data["retrieved_images"]:
46             pil_img = to_pil_image(img)
47             if pil_img:
48                 pil_img = pil_img.resize((128, 128), Image.LANCZOS)
49                 retrieved_image_base64s.append(image_to_base64_data_url(pil_img))
50                 print("Retrieved image converted successfully")
51
52     sections.append(section("Retrieved Images",
53                             f"{len(retrieved_image_base64s)} image(s) attached." if
54                             retrieved_image_base64s else "<no usable retrieved images>"))
55
56     # --- Retrieved Audios ---
57     if llm_data.get("retrieved_audios"):
58         raud_list = "\n\n".join(
59             [f"Retrieved Audio {i+1} (Transcription):\n{t}" for i, t in
60              enumerate(llm_data["retrieved_audios"])]
61         )
62         sections.append(section("Retrieved Audio Transcriptions", raud_list))
63
64     # --- Merge all context ---
65     full_context = "\n\n".join(sections).strip() or "No multimodal content provided."
66
67     # === Construct message content ===
68     content_list = [{"type": "text", "text": full_context}]
69
70     # Add images as base64 URLs
71     for img_url in image_base64s + retrieved_image_base64s:
72         content_list.append({
73             "type": "image_url",
74             "image_url": {"url": img_url}
75         })
76
77     # === Call the LLM ---
78     try:
79         completion = client.chat.completions.create(
80             model=LLM_MODEL_NAME, # e.g. "openai/gpt-5-image-mini"
81             extra_headers={
82                 "HTTP-Referer": "http://localhost",
83                 "X-Title": "Multimodal RAG Chatbot"
84             },
85             messages=[
86                 {
87                     "role": "system",
88                     "content": (
89                         "You are a multimodal reasoning assistant.\n"
90                         "You receive user inputs and retrieved multimodal data (text,
91                         image, audio).\n"
92                         "If the user asked a question, answer it using relevant data.\n"
93                         "If not, describe the multimodal inputs clearly in Markdown."
94                     )
95                 },
96                 {
97                     "role": "user",
98                     "content": content_list
99                 }
100             ]
101         )

```

```

102     response_text = completion.choices[0].message.content
103     except Exception as e:
104         response_text = f"LLM call failed: {e}"
105
106     # === Logging ===
107     log_dir = os.path.join(PROJECT_ROOT_PATH, "log")
108     os.makedirs(log_dir, exist_ok=True)
109     log_path = os.path.join(log_dir, "llm_full_context.txt")
110     try:
111         with open(log_path, "w", encoding="utf-8") as f:
112             f.write("=== FULL CONTEXT SENT TO LLM ===\n\n")
113             f.write(full_context + "\n\n")
114             f.write("=" * 50 + "\n\n")
115             f.write("=== LLM RESPONSE ===\n\n")
116             f.write(response_text + "\n")
117         print(f"Full context written to {log_path}")
118     except Exception as e:
119         print(f"Failed to write LLM log: {e}")
120
121     return response_text

```

این تابع `feed_data_into_llm()` یکی از مؤلفه‌های اصلی در معماری یک سیستم *Multimodal RAG Chatbot* محسوب می‌شود. وظیفه آن تجمیع، آماده‌سازی و ارسال تمامی داده‌های چندمدالی (شامل متن، تصویر، و صوت) به مدل زبان بزرگ (*LLM*) و بازگرداندن پاسخ نهایی مدل است.

تعریف تابع: `def feed_data_into_llm(llm_data: dict) -> str:`

- ورودی: یک دیکشنری (*dictionary*) به نام `llm_data` که حاوی داده‌های چندمدالی پروژه است (این دیکشنری در تابع `search_multimodal` تولید می‌شود).
- خروجی: یک رشته (*string*) که شامل پاسخ نهایی تولید شده توسط *LLM* است.

تابع کمکی داخلی: `def section(title: str, content: str) -> str: section()`

- این تابع به منظور قالب‌بندی محتوا در فرمت *Markdown* طراحی شده است.
- از آن برای تفکیک و ارائه منظم ورودی‌ها به *LLM* با استفاده از عنوان و محتوا استفاده می‌شود.

تشریح مراحل آماده‌سازی داده

۱. ورودی متنی کاربر (`user_text_query`): در صورتی که ورودی متنی از سوی کاربر دریافت شده باشد، این متن به عنوان یک مؤلفه مجزا در لیست *sections* ثبت می‌گردد.

```

1 if llm_data.get("user_text_query"):
2     sections.append(section("User Text Query", llm_data["user_text_query"]))
3

```

۲. تصاویر کاربر (`user_image_queries`):

- تصاویر ورودی به تابع `to_pil_image()` ارسال و به فرمت *PIL* تبدیل می‌گردند. تصاویر ورودی کاربر از طریق تابع `to_pil_image()` به فرمت استاندارد *PIL* (*Python Imaging Library*) تبدیل می‌شوند.
- جهت کاهش حجم ورودی، تصاویر به ابعاد 128×128 تغییر اندازه داده می‌شوند. برای بهینه‌سازی بار پردازشی و کاهش تأخیر، ابعاد تصاویر ورودی به صورت یکنواخت به 128×128 پیکسل (با استفاده از متد `Image.LANCZOS`) تغییر اندازه می‌یابد.
- تصاویر با استفاده از `image_to_base64_data_url()` به رشته‌ی *Base64* تبدیل می‌شوند تا قابلیت ارسال در بدنه *JSON* به مدل را داشته باشند. سپس، تصاویر به کد *Base64* تبدیل می‌شوند تا بتوان آن‌ها را در قالب *JSON* و مطابق با الزامات *API* مدل به *LLM* ارسال نمود.

```

1 image_base64s = []
2 if llm_data.get("user_image_queries"):
3     for img in llm_data["user_image_queries"]:
4         pil_img = to_pil_image(img)
5         if pil_img:
6             pil_img = pil_img.resize((128, 128), Image.LANCZOS)
7             image_base64s.append(image_to_base64_data_url(pil_img))

```



```

8 sections.append(section("User Image Queries", f"{len(image_base64s)} image(s)
9 attached." if image_base64s else "<no usable images>"))

```

۳. صوت کاربر (*user_audio_queries*): در صورت وجود ورودی صوتی، فرض بر این است که داده‌های صوتی قبلاً به متن (*Transcription*) تبدیل شده‌اند و ترنسکرپت‌های حاصله در قالب *Markdown* به منظور ورود به *LLM* آماده‌سازی می‌گردند.

```

1 if llm_data.get("user_audio_queries"):
2     aud_list = "\n\n".join(
3         [f"Audio {i+1} (Transcription):\n{t}" for i, t in
4          enumerate(llm_data["user_audio_queries"])]
5     )
6     sections.append(section("User Audio Queries", aud_list))

```

۴. متون بازیابی شده (*retrieved texts*) (*RAG*): متون مرتبطی که از طریق سیستم بازیابی اطلاعات (*Retrieval*) (مانند *Weaviate*) از پایگاه داده به دست آمده‌اند، به عنوان زمینه (*context*) جهت افزایش دقت و مرتبط بودن پاسخ *LLM*، به لیست بخش‌ها افزوده می‌گردند. این متون در واقع *context retrieval* محسوب می‌شوند.

```

1 if llm_data.get("retrieved_texts"):
2     txt_list = "\n\n".join(
3         [f"Retrieved Text {i+1}:\n{t}" for i, t in
4          enumerate(llm_data["retrieved_texts"])]
5     )
6     sections.append(section("Retrieved Texts", txt_list))

```

۵. تصاویر بازیابی شده (*retrieved images*): تصاویر مرتبط بازیابی شده از پایگاه داده، مشابه تصاویر کاربر، پس از تغییر اندازه و تبدیل به *Base64*، به ورودی نهایی مدل اضافه می‌شوند.

```

1 retrieved_image_base64s = []
2 for img in llm_data.get("retrieved_images", []):
3     pil_img = to_pil_image(img)
4     if pil_img:
5         pil_img = pil_img.resize((128, 128), Image.LANCZOS)
6         retrieved_image_base64s.append(image_to_base64_data_url(pil_img))
7 sections.append(section("Retrieved Images", f"{len(retrieved_image_base64s)}
8 image(s) attached." if retrieved_image_base64s else "<no usable retrieved
images>"))

```

۶. صوت‌های بازیابی شده (*retrieved audios*): ترنسکرپت‌های صوتی بازیابی شده به متن نهایی افزوده می‌شوند تا *LLM* بتواند از اطلاعات شنیداری مرتبط با جستجو نیز استفاده نماید.

```

1 if llm_data.get("retrieved_audios"):
2     raud_list = "\n\n".join(
3         [f"Retrieved Audio {i+1} (Transcription):\n{t}" for i, t in
4          enumerate(llm_data["retrieved_audios"])]
5     )
6     sections.append(section("Retrieved Audio Transcriptions", raud_list))

```

ترکیب داده‌ها و فراخوانی *LLM*

- ترکیب زمینه: تمامی بخش‌های آماده‌سازی شده از ورودی‌های کاربر و داده‌های بازیابی شده به یک رشته واحد (*full_context*) تلفیق می‌گردند تا به عنوان ورودی متنی جامع به مدل ارسال شوند.

```

1 full_context = "\n\n".join(sections).strip() or "No multimodal content provided."
2

```

- ساخت پیام: لیست پیام‌های ارسالی به *LLM* شامل *full_context* به عنوان متن و رشته‌های *Base64* تصاویر (شامل تصاویر کاربر و تصاویر بازیابی شده) است.

```

1 content_list = [{"type": "text", "text": full_context}]
2 for img_url in image_base64s + retrieved_image_base64s:
3     content_list.append({"type": "image_url", "image_url": {"url": img_url}})
4

```

- فراخوانی مدل: درخواست به مدل زبانی از طریق تابع `client.chat.completions.create` ارسال می‌شود و شامل دو پیام اصلی است: پیام `system` (جهت تعریف نقش دستیار) و پیام `user` (شامل `content_list` که حاوی داده‌های چندمددالی است).

```

1 completion = client.chat.completions.create(
2     model=LLM_MODEL_NAME,
3     extra_headers={
4         "HTTP-Referer": "http://localhost",
5         "X-Title": "Multimodal RAG Chatbot"
6     },
7     messages=[
8         {
9             "role": "system",
10            "content": (
11                "You are a multimodal reasoning assistant.\n"
12                "You receive user inputs and retrieved multimodal data (text,
13                image, audio).\n"
14                "If the user asked a question, answer it using relevant data.\n"
15                "If not, describe the multimodal inputs clearly in Markdown."
16            )
17        },
18        {
19            "role": "user",
20            "content": content_list
21        }
22    ]
23 )

```

- پاسخ: مدل بر اساس داده‌های ورودی، پاسخ نهایی را تولید می‌کند: `response_text = completion.choices[0].message.content`.

مکانیسم ثبت لاگ (*Logging*) کل ورودی (`full_context`) و خروجی (`response_text`)، به منظور اشکال‌زدایی و مستندسازی جریان داده، در فایل `llm_full_context.txt` واقع در پوشه `log` ذخیره می‌گردند.

```

1 log_path = os.path.join(PROJECT_ROOT_PATH, "log", "llm_full_context.txt")

```

نتیجه نهایی تابع، پاسخ نهایی تولیدشده توسط `LLM` (`response_text`) را بازمی‌گرداند تا جهت نمایش در واسط کاربری (*UI*) یا محیط چت‌بات (*Chatbot*) مورد استفاده قرار گیرد.

```

1 return response_text

```

۳.۶.۷ توابع کمکی

در تابع `feed_data_into_llm` از توابع زیر استفاده شده است:
تابع `to_pil_image`

```

1 def to_pil_image(img_candidate):
2     """Convert input to a PIL.Image (RGB)."""
3     if img_candidate is None:
4         return None
5     if isinstance(img_candidate, Image.Image):
6         return img_candidate.convert("RGB")
7     if isinstance(img_candidate, (bytes, bytearray)):
8         try:
9             return Image.open(io.BytesIO(img_candidate)).convert("RGB")
10        except Exception:
11            return None
12     if isinstance(img_candidate, str):
13         if os.path.exists(img_candidate):
14             try:

```

```

15         return Image.open(img_candidate).convert("RGB")
16     except Exception:
17         return None
18     if img_candidate.startswith("http://") or img_candidate.startswith("https://"):
19         try:
20             resp = requests.get(img_candidate, timeout=5)
21             if resp.status_code == 200:
22                 return Image.open(io.BytesIO(resp.content)).convert("RGB")
23             except Exception:
24                 return None
25         return None
26     if isinstance(img_candidate, dict):
27         b = img_candidate.get("bytes") or img_candidate.get("content") or
img_candidate.get("data")
28         if isinstance(b, (bytes, bytearray)):
29             try:
30                 return Image.open(io.BytesIO(b)).convert("RGB")
31             except Exception:
32                 return None
33         path = img_candidate.get("path") or img_candidate.get("file")
34         if isinstance(path, str) and os.path.exists(path):
35             try:
36                 return Image.open(path).convert("RGB")
37             except Exception:
38                 return None
39         return None
40     return None

```

تابع `to_pil_image(img_candidate)` یک تبدیل‌کننده‌ی همه‌منظوره است که تلاش می‌کند هر ورودی تصویری ممکن (شامل فایل‌ها، بایت‌ها، URL‌ها، و داده‌های ساختاریافته) را به یک شیء تصویر از نوع `PIL.Image` (در حالت رنگی RGB) تبدیل نماید.

خلاصه عملکرد `to_pil_image` to مرحله به مرحله

حالت ورودی	عملکرد تابع
<code>None</code>	مقدار <code>None</code> بازگردانده می‌شود.
نوع <code>PIL.Image</code>	تصویر ورودی به حالت رنگی "RGB" تبدیل شده و بازگردانده می‌شود.
نوع <code>bytes</code> یا <code>bytearray</code>	تلاش می‌شود تصویر از داده‌ی باینری بارگذاری شود (<code>Image.open(io.BytesIO(...))</code>).
نوع <code>str</code> (رشته)	<ul style="list-style-type: none"> اگر مسیر فایل (<code>path</code>) باشد: تصویر از فایل محلی بارگذاری می‌شود. اگر <code>URL</code> باشد: تصویر از طریق شبکه دانلود و باز می‌شود.
نوع <code>dict</code>	داده‌های تصویر از کلیدهای "bytes"، "data"، "content"، یا "path" استخراج و شیء تصویر ساخته می‌شود.
سایر انواع	مقدار <code>None</code> بازگردانده می‌شود.

نتیجه: این تابع به عنوان یک رابط استاندارد ساز عمل می‌کند تا تمامی فرمت‌های ورودی تصویر برای پردازش‌های بعدی در سیستم، یکپارچه‌سازی شوند.
تابع `image_to_base64_data_url`

```

1 def image_to_base64_data_url(pil_img):
2     """Convert PIL.Image to base64 data URL (JPEG)."""
3     buffer = io.BytesIO()
4     pil_img.save(buffer, format="JPEG")
5     img_str = base64.b64encode(buffer.getvalue()).decode("utf-8")
6     return f"data:image/jpeg;base64,{img_str}"

```

این تابع تصویر ورودی (`PIL.Image`) را به یک رشته‌ی `Base64 Data URL` (`data:image/jpeg;base64,...`) تبدیل می‌کند که استاندارد الزامی برای ارسال داده‌های تصویری به مدل‌های چندمدالی از طریق `API` است.

تشریح مفهوم `Base64` `Base64` یک روش کدگذاری است که داده‌های باینری (نظیر تصاویر) را به رشته‌های متنی `ASCII` تبدیل می‌کند. این فرآیند حیاتی است زیرا مدل‌های زبان بزرگ (`LLM`) و اکثر `API`‌ها، داده‌ها را به طور عمده با فرمت متنی (`string`) پردازش و منتقل می‌نمایند.

لزوم تبدیل به Base64

- انتقال امن: *LLM* یا *API* توانایی پردازش مستقیم فایل‌های باینری (مانند *.jpg* یا *.png*) را در ساختار پیام‌های *JSON* ندارند.
- استانداردسازی: با تبدیل به *Base64*، تصویر به صورت یک رشته‌ی متنی ایمن و قابل انتقال در پروتکل‌های تحت وب در می‌آید.
- قابلیت بازیابی: در سمت مقصد (مدل یا سرور)، داده‌های *Base64* دوباره به محتوای باینری تصویر تبدیل و جهت پردازش استفاده می‌شوند.