

Matthew Xu and Vy Le

Lab 2 DATA266

2 May 2025

DATA 266 Lab 2

Part 1

MultiModal Retrieval-Augmented Generation

We propose two different approaches for building a MultiModal Retrieval-Augmented Generation (RAG) system that processes a dataset of economic reports, tables, and images to retrieve relevant data and generate meaningful responses to economic queries. The first architecture integrates PyPDF2, SentenceTransformer, scikit-learn's TF-IDF, ChromaDB vector database, and OpenAI's API to create a hybrid retrieval system that employs domain-specific query expansion and context-aware reranking for economic document analysis. The second method uses similar technologies, but employs LangChain's framework to create a comprehensive pipeline for processing PDF documents containing both textual and visual economical information.

The Embedding Pipeline

The embedding pipeline uses the SentenceTransformer library to generate dense vector representations of text, with a preference for the 'intfloat/e5-large-v2' model from Hugging Face. This model has demonstrated superior performance in semantic retrieval tasks compared to earlier models. The e5-large-v2 model, with approximately 335 million parameters, produces 1024-dimensional embeddings that effectively capture semantic relationships between economic concepts, effective on different terminology.

The embedding process is applied differently to various content types - text chunks receive straightforward embedding, while figures and tables undergo a specialized embedding approach that combines their captions with inferred descriptions. For example, a table is embedded not just with its caption but with an expanded description to enhance retrieval relevance.

```

# Initialize Chroma client
self.chroma_client = chromadb.PersistentClient(path=persist_directory)

# Create collections for different content types
self.text_collection = self.chroma_client.get_or_create_collection(
    name="text_chunks",
    embedding_function=embedding_functions.SentenceTransformerEmbeddingFunction(model_name='intfloat/e5-large-v2')
)

self.figure_collection = self.chroma_client.get_or_create_collection(
    name="figures",
    embedding_function=embedding_functions.SentenceTransformerEmbeddingFunction(model_name='intfloat/e5-large-v2')
)

self.table_collection = self.chroma_client.get_or_create_collection(
    name="tables",
    embedding_function=embedding_functions.SentenceTransformerEmbeddingFunction(model_name='intfloat/e5-large-v2')
)

```

To complement these dense vectors, the system also employs TF-IDF (Term Frequency-Inverse Document Frequency) vectorization from scikit-learn. This creates sparse vectors that excel at capturing exact matches and rare economic terminology that might get diluted in dense embeddings. We apply TF-IDF with a maximum of 512 features, focusing on the most discriminative terms in the corpus. This hybrid embedding approach combines the semantic understanding of transformers with the precision of TF-IDF—creates a more robust representation space particularly suited to economic documents where both conceptual understanding and terminological precision matter greatly.

```

def create_hybrid_embeddings(self):
    # Create sparse TF-IDF embeddings
    documents = [chunk['text'] for chunk in self.text_chunks]
    self.tfidf_vectorizer = TfidfVectorizer(max_features=512)
    self.tfidf_matrix = self.tfidf_vectorizer.fit_transform(documents)

    # Store text chunks in Chroma
    self.store_text_chunks_in_chroma()

    # Store figures and tables in Chroma
    self.store_figures_and_tables_in_chroma()

```

The Retrieval Architecture

The retrieval architecture in PrecisionRAG demonstrates a multi-stage approach optimized for economic question answering. The process begins with query expansion, where the original query is enriched with domain-specific terminology based on question classification. For instance, queries about financial crises are expanded with terms like "housing market" and "subprime mortgage," while queries about unemployment incorporate terms like "jobless rate"

and "economic hardship." This expansion is informed by domain expertise in economics and significantly improves recall by bridging vocabulary gaps.

```
def adjust_query_for_question(self, query, question_id):
    # Economic domain knowledge expansions
    expansion_mapping = {
        1: ["financial crisis", "housing market", "subprime mortgage", "lehman brothers"],
        2: ["unemployment impact", "jobless rate", "economic hardship"],
        3: ["real GDP", "nominal GDP", "constant prices", "price adjustment"],
        4: ["global recession", "world output", "advanced economies contraction"],
        5: ["unemployment rate increase", "job losses", "U.S. unemployment"],
        6: ["China growth rate", "Chinese economic expansion", "fiscal stimulus"],
        7: ["stock market crash", "equity prices", "market decline"],
        8: ["Okun's law", "growth unemployment relationship", "output unemployment"],
        9: ["consumer price index", "inflation", "GDP deflator", "CPI"],
        10: ["European unemployment", "labor market rigidity", "Euro area"],
        11: ["China fiscal stimulus", "growth maintenance", "Chinese economy crisis"]
    }

    # Generate expanded query
    if question_id in expansion_mapping:
        expansion_terms = expansion_mapping[question_id]
        return query + " " + ".join(expansion_terms)

    return query
```

The expanded query then undergoes a hybrid retrieval process, querying both the Chroma vector database (using dense embeddings) and the TF-IDF sparse representations. Results from both methods are combined with a preference for dense results, but unique sparse results are preserved to ensure broad coverage. This hybrid approach helps capture both semantic relationships and exact terminology matches.

```
def hybrid_retrieve(self, query, question_id, top_k=8):
    # Get vector-based results from Chroma
    chroma_results = self.retrieve_from_chroma(query, question_id, top_k)

    # Get TF-IDF results
    tfidf_results = self.retrieve_from_tfidf(query, top_k)

    # Combine results with a preference for dense results but including unique sparse results
    dense_ids = [item['id'] for item in chroma_results]

    # Add unique sparse results
    for item in tfidf_results:
        if item['id'] not in dense_ids:
            chroma_results.append(item)

    # Sort by score
    chroma_results.sort(key=lambda x: x['score'], reverse=True)

    return chroma_results[:top_k]
```

Retrieved results then undergo a multi-dimensional scoring adjustment based on several factors. Question-specific keywords receive boosting, with matches to terms like "financial crisis" or "unemployment rate" amplifying scores for relevant questions. Content containing numerical data receives an additional boost—recognizing that in economics, statistical evidence is particularly valuable. The system also considers source types, preferentially ranking figures for visually-oriented questions and tables for data-heavy inquiries.

```
def retrieve_from_chroma(self, query, question_id, top_k=8):
    # Query adjustments based on question_id
    adjusted_query = self.adjust_query_for_question(query, question_id)

    # Query each collection
    text_results = self.text_collection.query(
        query_texts=[adjusted_query],
        n_results=top_k
    )

    figure_results = self.figure_collection.query(
        query_texts=[adjusted_query],
        n_results=min(top_k // 2, 3) # Limit figures to avoid overwhelming text
    )

    table_results = self.table_collection.query(
        query_texts=[adjusted_query],
        n_results=min(top_k // 2, 3) # Limit tables to avoid overwhelming text
    )

    # Combine and format results
    combined_results = []
```

The architecture also includes specialized mechanisms for figure and table retrieval. The code maintains a mapping between question types and relevant visualizations, recognizing that economic questions often have characteristic visual representations. For instance, unemployment questions typically benefit from unemployment rate charts, while GDP-related questions may need tables showing growth figures. When direct mapping isn't available, the system dynamically analyzes content references to identify the most relevant visual elements. This multi-stage retrieval architecture—query expansion, hybrid retrieval, context-aware reranking, and specialized visual element selection—creates a retrieval system particularly attuned to the complexities of economic documents.

The BERTScore result of 0.61794 for the RAG system indicates moderate success in generating answers that align semantically with reference responses in the economic domain. BERTScore measures semantic similarity using contextual embeddings, making it more nuanced than exact

match metrics like BLEU or ROUGE. This middle-range score is likely reflective of the system's hybrid approach - the neural embeddings successfully capture broad semantic relationships, while the TF-IDF components help with economic terminology.

LangChain Framework

We also implemented a multimodal RAG system specifically designed for economic document analysis, leveraging LangChain's framework to create a comprehensive pipeline for processing PDF documents containing both textual and visual information. The implementation uses PyPDFLoader for text extraction and includes alternate methods for image extraction.

```
def extract_images_from_pdf(pdf_path: str) -> List[str]:  
  
    try:  
        # Using pdf2image if available  
        from pdf2image import convert_from_path  
        print("Using pdf2image to extract images...")  
  
        # Convert PDF to list of PIL images |  
        images = convert_from_path(pdf_path)  
        img_base64_list = []  
  
        # Process each image  
        for i, img in enumerate(images):  
            # Convert to base64  
            buffered = io.BytesIO()  
            img.save(buffered, format="PNG")  
            img_base64 = base64.b64encode(buffered.getvalue()).decode("utf-8")  
            img_base64_list.append(img_base64)  
  
    return img_base64_list  
  
except ImportError:  
    print("pdf2image not available. Using fallback method...")  
  
return img_base64_list
```

The system employs a content processing approach that generates summaries for both textual and visual elements. For text, it uses GPT-4o to create concise summaries optimized for retrieval, while for images, it uses the same model with vision capabilities to generate descriptive summaries of charts, graphs, and other visual elements commonly found in economic documents. These summaries serve as the basis for vector embeddings created using OpenAI's embedding models.

```
def generate_text_summaries(texts: List[str], tables: Optional[List] = [], summarize_texts: bool=True, model: ChatOpenAI = GPT_4o):
    # Prompt
    prompt_text = """You are an assistant tasked with summarizing tables and text for retrieval. \
    These summaries will be embedded and used to retrieve the raw text or table elements. \
    Give a concise summary of the table or text that is well optimized for retrieval. Table or text: {element} """
    prompt = PromptTemplate.from_template(prompt_text)
    empty_response = RunnableLambda(
        lambda x: AIMessage(content="Error processing document")
    )

    # Text summary chain
    summarize_chain = (
        {
            "element": lambda x: x
        }
        | prompt
        | model
        | StrOutputParser()
    )

    # Initialize empty summaries
    text_summaries = []
    table_summaries = []

    # Apply to text if texts are provided and summarization is requested
    if texts and summarize_texts:
        text_summaries = summarize_chain.batch(texts, {"max_concurrency": 1})
    elif texts:
        text_summaries = texts

    # Apply to tables if tables are provided
    if tables:
        table_summaries = summarize_chain.batch(tables, {"max_concurrency": 1})
        return text_summaries, table_summaries
    else:
        return text_summaries
```

```

def image_summarize(img_base64, prompt, model: ChatOpenAI = GPT_4o):
    msg = model.invoke(
        [
            HumanMessage(
                content=[
                    {"type": "text", "text": prompt},
                    {
                        "type": "image_url",
                        "image_url": {
                            "url": f"data:image/jpeg;base64,{img_base64}"
                        },
                    },
                ]
            )
        ]
    )
    return msg.content

def generate_img_summaries(img_base64_list):
    # Store image summaries
    image_summaries = []

    # Prompt
    prompt = """You are an assistant tasked with summarizing images for retrieval. \
    These summaries will be embedded and used to retrieve the raw image. \
    Give a concise summary of the image that is well optimized for retrieval."""

    # Apply to images
    for base64_image in img_base64_list:
        image_summaries.append(image_summarize(base64_image, prompt))

    return image_summaries

```

For knowledge storage and retrieval, the system implements a MultiVectorRetriever backed by ChromaDB. This approach separates the storage of summary embeddings from the original content, allowing the system to retrieve relevant documents based on semantic similarity of summaries while returning the original, unmodified content. The retrieval mechanism handles both text and images, determining which modality is most relevant to a given query. The RAG chain implementation demonstrates prompt engineering, positioning GPT-4o as a "distinguished macroeconomist" to enhance domain-specific responses. The chain dynamically constructs multimodal prompts that include both retrieved text and relevant images in base64 format, enabling the model to interpret charts and graphs when answering economic questions.

```

def create_multi_vector_retriever(vectorstore, text_summaries: List[str], texts: List[str], image_summaries, images, table_summaries: Optional[List[str]] = None, tables: Optional[List[str]] = None):
    # Initialize the storage layer
    store = InMemoryStore()
    id_key = "doc_id"

    # Create the multi-vector retriever
    retriever = MultiVectorRetriever(
        vectorstore=vectorstore,
        docstore=store,
        id_key=id_key,
    )

    # Check that text_summaries is not empty before adding
    if text_summaries:
        update_documents(retriever, text_summaries, texts)
    # Check that table_summaries is not empty before adding
    if table_summaries:
        update_documents(retriever, table_summaries, tables)
    # Check that image_summaries is not empty before adding
    if image_summaries:
        update_documents(retriever, image_summaries, images)

    return retriever

```

```

def img_prompt_func(data_dict):

    formatted_texts = "\n".join(data_dict["context"]["texts"])
    messages = []

    # Adding the text for analysis
    text_message = {
        "type": "text",
        "text": (
            "You are a distinguished macroeconomist guiding students in answering questions\n"
            "related to their macroeconomics course.\n"
            "You will be given a mixed of text, tables, and image(s) usually of charts or graphs.\n"
            "Use this information to answer questions related to their macroeconomics lecture and topics related to the course. \n"
            "Provide concise and to-the-point answers.\n"
            f"User-provided question: {data_dict['question']}\n\n"
            "Text and / or tables:\n"
            f"{formatted_texts}"
        ),
    }
    messages.append(text_message)
    # Adding image(s) to the messages if present
    if data_dict["context"]["images"]:
        for image in data_dict["context"]["images"]:
            image_message = {
                "type": "image_url",
                "image_url": {"url": f"data:image/jpeg;base64,{image}"}
            }
            messages.append(image_message)
    return [HumanMessage(content=messages)]

```

This implementation differs by using numerous helper functions for tasks like base64 encoding/decoding, image resizing, and type detection, creating a robust system capable of handling the complexities of mixed-modality content typical in economic literature and educational materials.

Conclusion

Comparing the two models, the Langchain RAG implementation achieved a higher BERTScore of 0.69341 compared to the previous RAG's 0.61794. This performance difference likely stems from distinct technical approaches: The previous RAG relies on hybrid embeddings combining SentenceTransformer dense vectors with TF-IDF sparse representations specifically optimized for economic text, while the multimodal Langchain system leverages OpenAI's embedding models and GPT-4o's vision capabilities to process both textual content and visual elements like charts and graphs. The Langchain approach's ability to interpret visual economic data alongside

text appears to provide a meaningful advantage in addressing economic questions that often depend on data visualization for complete understanding.

Part 2

Fine-Tuning Stable Diffusion for Image Generation

Part 2 implements a fine-tuning pipeline for Stable Diffusion focused on generating high-quality food images. The pipeline uses the Food101 dataset with LoRA (Low-Rank Adaptation) to efficiently adapt the model while preserving the base model's capabilities. The implementation firstly creates a custom EnhancedFood101Dataset class that loads the Food101 dataset, creates varied text prompts for each food category, processes images to the required resolution, and tokenized text prompts for the model.

```
class EnhancedFood101Dataset(Dataset):
    def __init__(self, dataset_split, tokenizer, resolution=512, max_train_samples=None):
        self.dataset = load_dataset("food101", split=dataset_split)
        if max_train_samples:
            self.dataset = self.dataset.select(range(max_train_samples))

        self.tokenizer = tokenizer
        self.resolution = resolution

        # Get the class names
        self.class_names = self.dataset.features["label"].names

        # Prompts
        self.label_to_prompt = {
            label: [
                f"A professional photograph of {label.replace('_', ' ')} on a plate, food photography",
                f"High-quality food photo of {label.replace('_', ' ')}), styled by a chef",
                f"Delicious {label.replace('_', ' ')} photographed for a cookbook"
            ] for label in self.class_names
        }

        # Transform for images
        self.transform = transforms.Compose([
            transforms.Resize((resolution, resolution)),
            transforms.ToTensor(),
            transforms.Normalize([-0.5], [0.5]) # Normalize to [-1, 1]
        ])

    def __len__(self):
        return len(self.dataset)

    def __getitem__(self, idx):
        item = self.dataset[idx]
        image = item["image"]
        label_id = item["label"]

        # Convert label id to text prompt
        label_name = self.class_names[label_id]
        prompt_list = self.label_to_prompt[label_name]
        prompt = random.choice(prompt_list)

        # Process image
        if not isinstance(image, Image.Image):
            image = Image.fromarray(image)

        image = self.transform(image)

        # Tokenize text
        text_inputs = self.tokenizer(
            prompt,
            padding="max_length",
            max_length=self.tokenizer.model_max_length,
            truncation=True,
            return_tensors="pt"
        )
```

For fine-tuning, we use LoRA adaptation which targets specific modules in both UNet and text encoder components with a configuration of rank=16, alpha=16, and dropout=0.05. The VAE component is frozen to focus training on the text-to-image mapping rather than general image compression.

```

def fine_tune_stable_diffusion():

    # Load base model
    model_id = "runwayml/stable-diffusion-v1-5"
    pipe = StableDiffusionPipeline.from_pretrained(model_id)

    # Fine-tune UNet and text encoder
    unet = pipe.unet
    text_encoder = pipe.text_encoder
    tokenizer = pipe.tokenizer
    vae = pipe.vae

    # Freeze VAE
    vae.requires_grad_(False)

    # Initialize noise scheduler for training
    noise_scheduler = DDPMscheduler.from_pretrained(model_id, subfolder="scheduler")

    # LoRA for UNet
    unet_lora_config = LoraConfig(
        r=16, # rank
        lora_alpha=16, # scaling
        target_modules=["to_q", "to_k", "to_v", "to_out.0"],
        lora_dropout=0.05,
        bias="none",
    )

    # LoRA for text encoder
    text_encoder_lora_config = LoraConfig(
        r=16,
        lora_alpha=16,
        target_modules=["q_proj", "k_proj", "v_proj", "out_proj"],
        lora_dropout=0.05,
        bias="none",
    )

    # LoRA configs
    unet = get_peft_model(unet, unet_lora_config)
    text_encoder = get_peft_model(text_encoder, text_encoder_lora_config)

```

The training process utilizes the AdamW optimizer with a learning rate of 5e-5, running for 5 epochs with a deliberate noise addition and prediction process in latent space. The base model chosen is "runwayml/stable-diffusion-v1-5" as it offers a good balance between quality and computational requirements. LoRA fine-tuning was selected because it dramatically reduces parameter count (training approximately 1% of parameters), preserves general capabilities while adapting to the food domain, enables faster training, and requires less storage for the resulting model.

```

# Train
unet.train()
text_encoder.train()

progress_bar = tqdm(range(num_epochs * len(train_dataloader)))
for epoch in range(num_epochs):
    for step, batch in enumerate(train_dataloader):
        # Get batch
        pixel_values = batch["pixel_values"].to(device, dtype=torch.float32)
        input_ids = batch["input_ids"].to(device)

        # Encode text
        with torch.no_grad():
            text_embeddings = text_encoder(input_ids)[0]

        # Encode image
        with torch.no_grad():
            # Get latent space representation
            latents = vae.encode(pixel_values).latent_dist.sample() * 0.18215

        # Add noise
        noise = torch.randn_like(latents)
        timesteps = torch.randint(0, noise_scheduler.config.num_train_timesteps,
                                  (latents.shape[0],), device=latents.device)
        noisy_latents = noise_scheduler.add_noise(latents, noise, timesteps)

        # Get prediction
        model_pred = unet(noisy_latents, timesteps, text_embeddings).sample

        # Loss
        loss = F.mse_loss(model_pred, noise, reduction="mean")

        # Update
        loss.backward()
        optimizer.step()
        optimizer.zero_grad()

        # Update progress
        progress_bar.update(1)
        if step % 10 == 0:
            print(f"Epoch: {epoch+1}, Step: {step}, Loss: {loss.item():.6f}")

# State dictionaries
unet_state_dict = unet.state_dict()
text_encoder_state_dict = text_encoder.state_dict()

```

For evaluation, CLIP Score was used to measure alignment between generated images and their prompts and Inception Score that assesses diversity and quality of generated images using a pretrained Inception v3 model.

The image generation uses diverse prompt engineering with 10 different style variations and varying generation parameters (guidance scales 7.0-9.0, step counts 40-50, different random seeds) to maximize diversity. It uses the optimized DPM++ Karras scheduler, to produce higher quality results than default schedulers and better maintains fine details in complex food textures.

```
# Generate images and evaluate
def generate_and_evaluate(pipe, unet, text_encoder, tokenizer, food_classes):

    # SD pipeline
    fine_tuned_pipe = StableDiffusionPipeline.from_pretrained(
        "runwayml/stable-diffusion-v1-5",
        unet=unet,
        text_encoder=text_encoder,
        tokenizer=tokenizer,
    )
    fine_tuned_pipe = fine_tuned_pipe.to(device)

    # DPM++ 2M Karras scheduler
    fine_tuned_pipe.scheduler = DPMSolverMultistepScheduler.from_config(
        fine_tuned_pipe.scheduler.config,
        algorithm_type="dpmsolver++",
        use_karras_sigmas=True
    )
```

```

print("Generating images...")
images = []
for i, prompt in enumerate(prompts):
    print(f"Generating image {i+1}/{len(prompts)}")

    image = fine_tuned_pipe(
        prompt,
        num_inference_steps=50,
        guidance_scale=8.0,
        negative_prompt="blurry, low quality, distorted, ugly, bad proportions, amateur"
    ).images[0]

    # Save
    image_path = f"{GENERATED_DIR}/image_{i+1}.png"
    image.save(image_path)
    images.append(image)

num_is_images = min(50, len(food_classes))
all_selected_foods = random.sample(food_classes, num_is_images)
is_prompts = []

# Style variation prompts
style_variations = [
    "professional food photography with studio lighting",
    "overhead view on rustic wooden table",
    "minimalist plating on white background",
    "gourmet restaurant presentation with garnish",
    "homemade style with natural lighting",
    "close-up macro photography showing texture",
    "artistic food styling with decorative elements",
    "street food style on paper wrapper",
    "high-end magazine editorial shot",
    "casual dining atmosphere"
]

# Prompt generation with high variation
for i, food in enumerate(all_selected_foods):
    style = style_variations[i % len(style_variations)]
    # Vary camera angles, lighting, and styling
    is_prompts.append(
        f"A {style}, showing {food.replace('_', ' ')}, appetizing, detailed, 8k"
    )

is_images = []
# Change generation parameters for each image to increase diversity
for i, prompt in enumerate(is_prompts):
    print(f"Generating IS image {i+1}/{len(is_prompts)}")

    # Vary the guidance scale to increase diversity
    guidance = 7.0 + (i % 3) # Values between 7.0 and 9.0

    # Vary the number of steps for different detail levels
    steps = 40 + (i % 3) * 5 # Values between 40 and 50

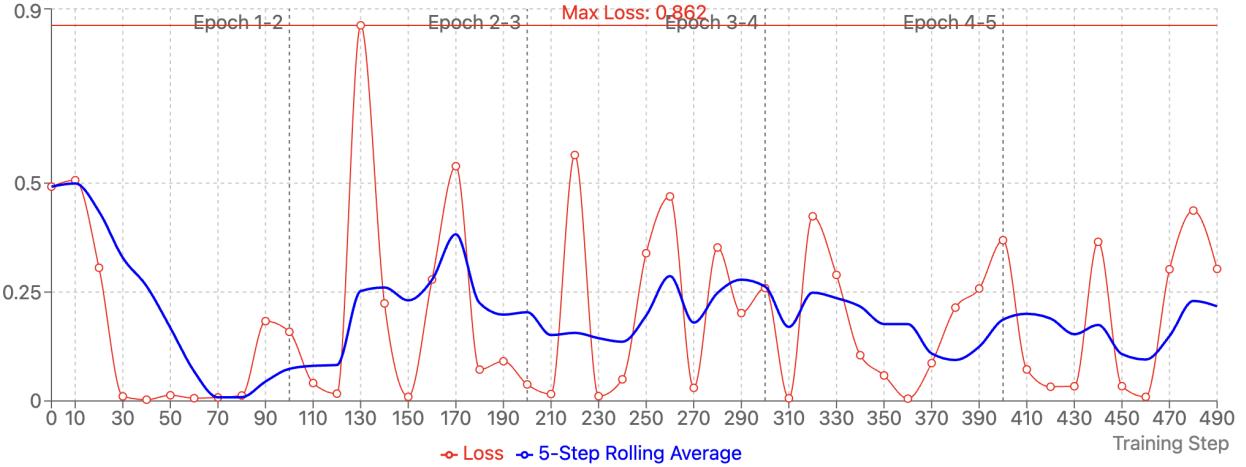
    # Use different seeds for each generation to maximize diversity
    generator = torch.Generator(device=device).manual_seed(i * 42 + 1000)

    image = fine_tuned_pipe(
        prompt,
        num_inference_steps=steps,
        guidance_scale=guidance,
        generator=generator,
        negative_prompt="blurry, ugly, deformed, noisy, low quality, oversaturated, high contrast"
    ).images[0]

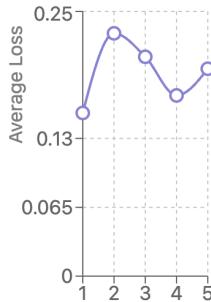
```

The Food101 dataset was chosen because it's a standard benchmark containing 101 food categories with 1,000 images per category, offering real-world food images with natural variation. The code saves fine-tuned model weights, LoRA configuration files, combined checkpoints with timestamps, generated sample images, evaluation metrics, and visualization plots.

Training Loss During Fine-Tuning



Epoch Average Losses



The early rapid loss convergence followed by occasional spikes is common in diffusion models, especially when using LoRA which can adapt rapidly to specific data patterns. The volatility in the later epochs could indicate that the model was encountering highly varied images in the Food101. This is actually a good sign, as it shows the model wasn't settling into a local minimum but continuing to learn from diverse examples and is typical in diffusion model training, as the loss depends greatly on the specific noise levels and images in each batch. The final loss of 0.303169 is relatively high compared to the minimum observed loss of 0.002988, which could mean the learning process was still ongoing and might have benefited from additional epochs.

Average CLIP similarity score: 32.3829
Reference points – Same content: 32.1390, Unrelated: 11.0407
Inception Score: 7.0800 ± 1.8223
Visualization saved to /content/drive/MyDrive/DATA266_Lab_2/food101_model/results.png
Inception Score visualization saved to /content/drive/MyDrive/DATA266_Lab_2/food101_model/inception_score_vis.png
Saved combined checkpoint to /content/drive/MyDrive/DATA266_Lab_2/food101_model/food101_sd_checkpoint_20250424_001019.pt
Fine-tuning and evaluation complete!

Prompt: A close-up, professional photo... Inception Score: 7.0800 ± 1.8223



CLIP Score: 32.1390

Prompt: A close-up, professional photo...



CLIP Score: 32.8308

Inception Score Visualization: 7.0800 ± 1.8223



The CLIP similarity score of 32.3829 indicates excellent alignment between generated images and their text prompts. For context, this score is slightly higher than the dataset (32.1390), suggesting the model generates images that match their text descriptions as well as or better. The dataset reference point is generated using a paired index from the food101 dataset itself, to provide context and standard to where the CLIP scores should be. The substantial gap between this score and the unrelated-content score (11.0407) confirms the model is creating images specifically referencing their respective prompts rather than generic food imagery.

The Inception Score of 7.0800 ± 1.8223 demonstrates good diversity and quality in the generated images. This metric measures both the clarity/quality of individual images and the diversity across the 50 generated images. The score indicates the model is producing varied, distinct food images rather than repetitive or similar outputs. The standard deviation (± 1.8223) suggests some variability in generation quality, which is expected when creating diverse food categories.

Part 3

Developing an Agentic AI Travel Assistant

The purpose of this part is to build up an Agentic AI which is known as the multi agent AI planner including flights, hotel and weather management. By using CrewAI and AI model, we aim to create a flight AI agent, weather AI agent and hotel AI agent separately to retrieve data from APIs platforms and create a centralized workflow to communicate these agents to produce an experienced and personalized travel itinerary based on user queries.

In short, there is a structure to perform the four agents:

```
Agentic-ai-travel/
    └── modules/
        ├── flight_main.py      # flight_agent + flight_task
        ├── hotel_main.py      # hotel_agent + hotel_task
        ├── weather_main.py    # weather_agent + weather_task
        └── planner_agent_main.py # planner_agent + planner_task (coordinates others)

    └── tools/
        ├── flight_tool.py     # API interaction logic for flights
        ├── hotel_tool.py      # API interaction logic for hotels
        └── weather_tool.py    # API interaction logic for weather

    └── configs/
        └── amadues_api.py     # API keys
```



README.md

Project overview and setup

Flight API Agent:

This agent will get the data from Amadeus API which is known as the API free platform to connect developers with the flight offer search API over 400 airlines, including detailed information such as airlines, departure, arrival, duration time and pricing. To connect with Amadeus API, I create an authorization like client ID and client secret to generate the access token via the url: <https://test.api.amadeus.com/v1/security/oauth2/token>. The access token is only valid for 30 minutes, so it needs to request a new one periodically.

The primary function of the Flight Tool is to search for flights and retrieve comprehensive flight information based on user input. It accepts parameters such as origin, destination, departure date, return date, number of travelers, and whether the flight should be non-stop. These parameters are passed to url "<https://test.api.amadeus.com/v2/shopping/flight-offers>" to retrieve flight information. The tool queries the Amadeus API and returns a curated list of available flight options that match the specified travel criteria.

In the flight main file, I define the flight agent by calling the search flights tool by using crewAI agent. CrewAI is a framework used to facilitate collaboration between AI agents, allowing them to work together on complex tasks like content generation, research, and even automating workflows. It provides a structure for defining agents, assigning tasks, and orchestrating their interactions, essentially enabling AI to operate in a "crew" or team. CrewAI integrates with the AI model to allow each agent to think, plan, and execute tasks based on its role, tools, and goals. Without specific LLM defining, crewAI is going to use gpt 3.5 turbo as the default generation model.

To test the flight agent, I tried the query "Find the 2 nonstop flights from NYC to LON for 2 people, departing on 2025-05-05 and returning on 2025-05-10." and the expected output "List 3 good flight options and recommend the best one.". The AI agent is expected not only to retrieve relevant flight options from the Amadeus API but also to evaluate and recommend the most suitable flight for the user based on the query criteria.

Result:

```

==== FINAL RESULTS ====
Here are 3 good nonstop flight options from NYC to LON for 2 people, departing on 2025-05-05 and returning on 2025-05-10:

1. **Flight Option 1**
- Total Price: 1044.34 EUR
- Departure:
  - From: JFK at 2025-05-05, 23:00
  - To: LHR at 2025-05-06, 11:10
  - Carrier: VS, Flight: 138
- Return:
  - From: LHR at 2025-05-10, 18:15
  - To: JFK at 2025-05-10, 21:00
  - Carrier: VS, Flight: 137
- Duration: 7h 10m (departure) + 7h 45m (return)

2. **Flight Option 2**
- Total Price: 1044.34 EUR
- Departure:
  - From: JFK at 2025-05-05, 21:00
  - To: LHR at 2025-05-06, 09:15
  - Carrier: VS, Flight: 46
- Return:
  - From: LHR at 2025-05-10, 18:15
  - To: JFK at 2025-05-10, 21:00
  - Carrier: VS, Flight: 137
- Duration: 7h 15m (departure) + 7h 45m (return)

3. **Flight Option 3**
- Total Price: 1044.34 EUR
- Departure:
  - From: JFK at 2025-05-05, 23:00
  - To: LHR at 2025-05-06, 11:10
  - Carrier: VS, Flight: 138
- Return:
  - From: LHR at 2025-05-10, 16:50
  - To: JFK at 2025-05-10, 19:40
  - Carrier: VS, Flight: 153
- Duration: 7h 10m (departure) + 7h 50m (return)

```

Among these options, ****Flight Option 1**** is the best choice due to its ideal departure and arrival times,

Weather API Agent:

Weather data is obtained from OpenWeatherMap API which is an online service and provides a current real-time weather and forecast. By using a free API service, there is limited access to retrieve weather forecasts within the next 4 days from today. It restricts selecting a certain or longer date range.

To access data from this API platform, an HTTP request is made by passing parameters into the URL. The required inputs include the API key, city name or location code, and the current date. The following URL can be used to retrieve temperature and weather conditions:

http://api.openweathermap.org/data/2.5/forecast?q={city}&appid={api_key}&units=metric. The get weather forecast tool will fetch weather data from the next 4 days from the given location and it will be called into the weather agent as the main task.

CrewAI acts as the weather agent, taking on the role of a weather assistant with the goal of forecasting weather conditions in a specified location for the next four days. The agent's

backstory establishes it as a reliable source for planning weather-dependent activities. A specific task was assigned to this agent: retrieve the weather forecast for London and assist in identifying the best day for outdoor plans.

Output:

```
==== FINAL RESULTS ====
The forecast for London is as follows:
- 2025-05-02: Temperature: 14.09°C, Condition: Overcast clouds
- 2025-05-03: Temperature: 10.58°C, Condition: Scattered clouds (Good day for going out)
- 2025-05-04: Temperature: 7.27°C, Condition: Broken clouds
- 2025-05-05: Temperature: 8.25°C, Condition: Scattered clouds
- 2025-05-06: Temperature: 9.8°C, Condition: Few clouds

Considering the weather conditions, May 3rd is the best day for going out, as it has scattered clouds and a mild temperature of 10.58°C.
```

Hotel API Agent:

Using the same Amadeus API as the Flight agent, the Hotel agent can explore over 150,000 accommodations worldwide. This API provides comprehensive information on hotel pricing, amenities, and user reviews based on the selected dates and location. Since both agents share the same platform, the Hotel agent can reuse the existing authentication method and access tokens, ensuring seamless integration and efficient API communication.

There are the following steps to retrieve hotel information:

The process begins by retrieving a list of hotel IDs for a given city using the endpoint: <https://test.api.amadeus.com/v1/reference-data/locations/hotels/by-city>, where the city code is passed as a parameter.

These hotel IDs are then used to fetch available hotel offers via the endpoint: <https://test.api.amadeus.com/v3/shopping/hotel-offers>, along with parameters such as check-in date, check-out date, and the number of adults.

Finally, user reviews and ratings for each hotel are retrieved using: <https://test.api.amadeus.com/v2/e-reputation/hotel-sentiments>. By combining data from these three endpoints, the system can present comprehensive hotel details including the hotel name, location, price, room type, description, and guest feedback.

CrewAI will define hotel agent as the role of hotel travel assistant with the goal is to help users find and compare hotel offers with ratings in a given city and date range. This hotel agent uses the tool 'get_hotel_data_combined' to perform the specific task like "Search for hotels in the city code 'LON' from 2025-05-05 to 2025-05-10 for 2 adults. Display name, location, price, room type, and rating.". The expectation output is set to list 3 good hotels and recommend the best one based on ratings, prices, and room descriptions.

Result:

```

==== FINAL RESULTS ====
1. **London Marriott Hotel Canary Wharf**
   - **Location:** LON | 51.50795, -0.02404
   - **Price:** **1438.30 GBP**
   - **Room Type:** Deluxe Room, 1 King, 31sqm/334sqft, Wireless internet (for a fee), Coffee/tea maker.

2. **London Marriott Marble Arch**
   - **Location:** LON | 51.51587, -0.16413
   - **Price:** **1316.25 GBP**
   - **Room Type:** Deluxe Room, 1 King, 22sqm/237sqft, Living/sitting area, Wireless internet (for a fee).

3. **London Marriott Hotel Park Lane**
   - **Location:** LON | 51.51296, -0.15772
   - **Price:** **2133.50 GBP**
   - **Room Type:** Deluxe Room, 1 King, 28sqm/301sqft, Wireless internet (for a fee), Coffee/tea maker.

4. **London Marriott Hotel Regents Park** *(Recommended)*
   - **Location:** LON | 51.54187, -0.17099
   - **Price:** **1160.25 GBP**
   - **Room Type:** Deluxe room, balcony, 1 King, 21sqm/226sqft-28sqm/301sqft, Wireless internet.

The best recommendation based on ratings, prices, and room descriptions is the **London Marriott Hotel Regents Park** for **1160.25 GBP**.

```

Planner API Agent:

After developing three functional AI systems, it is required to build a central workflow to facilitate the communication between these agents to deliver a good itinerary plan for travel. Therefore, we create a planner travel agent to manage the plan by reusing the tools from these agents like search flights, get weather data and search hotels.

The main planner agent is defined as the travel planner assistant, the goal is to create a full travel plan with flights, hotels, and weather forecast. This will call tools: `search_flights` from flight tool, `get_hotel_data_combined` from hotel tool and `get_weather_data` from weather tool. With the user query:

Plan a complete trip from SFO to PAR from 2025-05-05 to 2025-05-10 for 2 people.

Use tools to find:

1. Flight options.
2. Best hotel options.
3. Weather forecasts.

An expected plan is a full travel plan itinerary with recommendations for the best 2 flights, hotels, and weather.

CrewAI will start the process to generate the travel plan from flights, then hotel and weather based on the tools provided. The parameters are automatically detected from the query like origin (San Francisco), destination (Paris), departure date/ check-in (2025-05-05), return date/ check-out (2025-05-10), 2 adults. Then it combines these information together as the final result and the AI model suggests the option based on the plan.

Result:

==== COMPLETE TRAVEL PLAN ===

Travel Itinerary: San Francisco (SFO) to Paris (PAR)

Flight Options:

Option 1:

- **Departure:**

- From: SFO at 2025-05-05, 10:16 AM
- To: EWR at 2025-05-05, 6:54 PM
- Carrier: AS, Flight Number: 219
- From: EWR at 2025-05-05, 11:00 PM
- To: ORY at 2025-05-06, 12:15 PM
- Carrier: BF, Flight Number: 721

- **Total Price:** 1486.48 EUR

- **Duration:** 16 hours 59 minutes

Option 2:

- **Departure:**

- From: SFO at 2025-05-05, 6:00 AM
- To: EWR at 2025-05-05, 2:29 PM
- Carrier: AS, Flight Number: 292
- From: EWR at 2025-05-05, 11:00 PM
- To: ORY at 2025-05-06, 12:15 PM
- Carrier: BF, Flight Number: 721

- **Total Price:** 1486.48 EUR

- **Duration:** 21 hours 15 minutes

Hotel Options in Paris:

1. Renaissance Paris Vendome Hotel

- **Check-in:** 2025-05-05
- **Check-out:** 2025-05-10
- **Price:** 2702.05 EUR
- **Room Detail:** AP7 - Prepay Non-refundable, Mini fridge, 18sqm

2. Best Western Gaillon Opera

- **Check-in:** 2025-05-05
- **Check-out:** 2025-05-10
- **Price:** 811.50 EUR
- **Room Detail:** 1 DOUBLE BED, NSMK, STND

Weather Forecast for Paris (May 5-10, 2025):

- **May 5:** Light Rain, Temp: 10.74°C
- **May 6:** Broken Clouds, Temp: 12.44°C
- **May 7:** Condition: TBA (not available)
- **May 8:** Condition: TBA (not available)
- **May 9:** Condition: TBA (not available)
- **May 10:** Condition: TBA (not available)

Note: Weather for additional days may not be fully available yet.

Recommendations:

1. **Flight Option:** Choose either of the flight options to minimize layover duration.
2. **Hotel Option:** Consider the Best Western Gaillon Opera for a cost-effective choice or the Renaissance Paris Vendome Hotel
3. **Pack accordingly:** Given the light rain forecast, bring an umbrella and waterproof clothing.