# AIM825 Course Project: Multimodal Visual Question Answering with Amazon Berkeley Objects Dataset

Shreyas S          Tanay Nagarkar          R Lakshman
IMT2022078          IMT2022083          IMT2022090

May 18, 2025

## 1  Overview

This project implements a Visual Question Answering (VQA) system using two approaches:

- Primary Model: BLIP with LoRA fine-tuning. Bootstrapped Language-Image Pretraining (BLIP) is a model designed for **vision-language tasks**, including **Visual Question Answering (VQA)**, **image captioning**, and **image-text retrieval**. BLIP leverages **contrastive learning**, **captioning**, and **question answering** to improve multimodal understanding.

- Experimentational model : CLIP for semantic similarity scoring

  Contrastive Language-Image Pretraining (CLIP) is a **discriminative model** designed for **image-text matching** tasks. Unlike generative models, CLIP does not generate textual responses but instead computes **similarity scores** between images and text embeddings. CLIP is particularly effective for **semantic similarity-based retrieval**, where it ranks the most relevant textual descriptions corresponding to a given image.

  CLIP has been widely used in tasks such as **zero-shot image classification**, **image retrieval**, and **text-based search**, demonstrating strong capabilities in aligning vision and language representations.

We developed a custom dataset generation pipeline using multimodal prompts, implemented parameter-efficient fine-tuning with LoRA, and conducted comprehensive evaluation using both traditional and semantic metrics.

# 2 Data Curation

## 2.1 Prompt Engineering

We used the following structured prompt to generate diverse questions:

```
prompt = ("For the given image, generate **exactly** 5
    diverse questions that can be answered just by looking at
    it. "
        "Ensure a mix of simple and difficult questions.
    Provide **only** the question followed by a single-word
    answer, separated by a colon.")
```

This prompt ensured:

- Fixed number of questions per image (5)

- Balanced difficulty levels

- Strict formatting for automated parsing

- Single-word answers for simplified evaluation

## 2.2 Dataset Processing Pipeline

The data processing involved:

```
# BLIP Data Loading Implementation
df = pd.read_csv(CSV_PATH)
with open(JSON_PATH, "r") as f:
    qna_dict = json.load(f)

vqa_data = []
for _, row in df.iterrows():
    image_id = str(row["image_id"])
    image_path = os.path.join(IMAGE_ROOT, row["path"])
    if not os.path.exists(image_path): continue
    for qa in qna_dict.get(image_id, []):
        q = qa["question"].strip("?").strip() + "?"
```

```
13          a = qa["answer"].strip()
14          if a:
15              vqa_data.append({
16                  "image_path": image_path,
17                  "question": q,
18                  "answer": a
19              })
```

Key processing steps:

- Image path validation and cleaning

- Question standardization (ensuring proper punctuation)

- Answer normalization (strip whitespace)

- Type conversion for image IDs

## 2.3    Dataset Statistics

We firstly tried a small portion of the dataset for our BLIP model. We generated **2,500 questions** for the BLIP model using Gemini. Out of these, **2,000 questions** were used for training, and **500 questions** were used for testing. Further evaluation was conducted based on the model's performance on the test set.

# 3    Model Architectures

## 3.1    BLIP-VQA Implementation

The BLIP model was implemented with the following key components:

```
1  # BLIP Model Initialization
2  processor = BlipProcessor.from_pretrained(MODEL_NAME)
3  model = BlipForQuestionAnswering.from_pretrained(MODEL_NAME).
       to(device)
4
5  # LoRA Configuration
6  lora_config = LoraConfig(
7      r=8,
8      lora_alpha=16,
9      target_modules=["query", "value"],
10     lora_dropout=0.05,
```

```
11      bias="none",
12  )
13
14  model = get_peft_model(prepare_model_for_kbit_training(model)
        , lora_config)
```

Key architectural decisions:

- Used base BLIP-VQA model (Salesforce/blip-vqa-base)

- Applied LoRA only to query and value projections

- Maintained original model dimensionality (768 hidden size)

- Kept vision encoder frozen during fine-tuning

## 3.2 CLIP Implementation

The CLIP implementation focused on semantic similarity:

```
1  # CLIP Feature Extraction
2  def encode_text_embeddings(question, answer):
3      text_input = processor(text=f"{question} {answer}",
4                          return_tensors="pt",
5                          padding=True).to(device)
6      with torch.no_grad():
7          return model.get_text_features(**text_input)
8
9  def load_image(image_id):
10     image_path = os.path.join(image_folder, image_id_to_path[
    image_id])
11     image = Image.open(image_path).convert("RGB")
12     return processor(images=image, return_tensors="pt").to(
    device)
```

# 4 Training Methodology

## 4.1 BLIP Fine-Tuning

The training process involved:

```
1  training_args = TrainingArguments(
2      output_dir="/content/blip_lora_finetuned",
3      per_device_train_batch_size=BATCH_SIZE,
```

```
4        per_device_eval_batch_size=BATCH_SIZE,
5        num_train_epochs=EPOCHS,
6        learning_rate=5e-5,
7        weight_decay=0.01,
8        warmup_steps=100,
9        lr_scheduler_type="cosine",
10       fp16=True,
11   )
12
13   trainer = Trainer(
14       model=model,
15       args=training_args,
16       train_dataset=train_dataset,
17       eval_dataset=eval_dataset
18   )
```

Key training parameters:

- Batch size: 4 (limited by GPU memory)

- Epochs: 5

- Learning rate: 5e-5 with cosine decay

- Mixed precision training (FP16)

- Weight decay for regularization

## 4.2   CLIP Training

The CLIP model was trained with:

```
1   optimizer = torch.optim.AdamW(model.parameters(), lr=5e-5)
2   criterion = torch.nn.BCEWithLogitsLoss()
3
4   for epoch in range(epochs):
5       for batch in train_loader:
6           optimizer.zero_grad()
7           image_embeds = model.get_image_features(pixel_values=
    batch["pixel_values"])
8           text_embeds = model.get_text_features(
9               input_ids=batch["input_ids"],
10              attention_mask=batch["attention_mask"]
11          )
12          loss = criterion((image_embeds * text_embeds).sum(dim
    =-1), batch["label"].float())
```

```
13          loss.backward()
14          optimizer.step()
```

# 5   Evaluation Metrics

In order to evaluate the performance of our VQA models, we employed three key metrics: **BERTScore F1**, **Exact Match**, and **BLEU-1**. Each metric was chosen based on its relevance to assessing generated answers in an image-based question-answering setup.

- **BERTScore F1**: Selected to measure **semantic similarity** between model-generated answers and reference answers. Unlike traditional n-gram matching metrics, BERTScore uses contextual embeddings to capture meaning rather than exact word matches, making it effective for VQA where wording variations can still convey the correct answer.

- **Exact Match**: Used as a strict evaluation metric, where an answer was considered correct only if it **exactly matched** the ground-truth reference answer. This metric is crucial for short-form answers, ensuring precision without relying on semantic closeness.

- **BLEU-1**: Chosen to assess the **n-gram overlap** between generated responses and reference answers. Since VQA answers are often short, BLEU evaluation was restricted to unigram precision (BLEU-1). This provided insights into how well the model reproduced expected words, despite single-word answers limiting higher-order n-gram effectiveness.

-  By using these three metrics in combination, we ensured a balanced evaluation covering **semantic correctness (BERTScore)**, **direct accuracy (Exact Match)**, and **word overlap (BLEU-1)**. This approach allowed for a comprehensive analysis of model performance across different aspects of answer quality.

## 5.1   Implementation Details

We implemented comprehensive evaluation metrics:

```
1 def evaluate_metrics(preds, refs):
2     bertscore = evaluate.load("bertscore")
```

```
3    bleu = evaluate.load("bleu")
4    rouge = evaluate.load("rouge")
5    meteor = evaluate.load("meteor")
6
7    metrics = {
8        "bertscore_f1": sum(bertscore.compute(
9            predictions=preds, references=refs, lang="en",
10           rescale_with_baseline=True)["f1"]) / len(preds),
11       "bleu": bleu.compute(predictions=preds, references=[[
    r] for r in refs])["bleu"],
12       "meteor": meteor.compute(predictions=preds,
    references=refs)["meteor"],
13       "rougeL": rouge.compute(predictions=preds, references
    =refs)["rougeL"]
14   }
15   return metrics
```
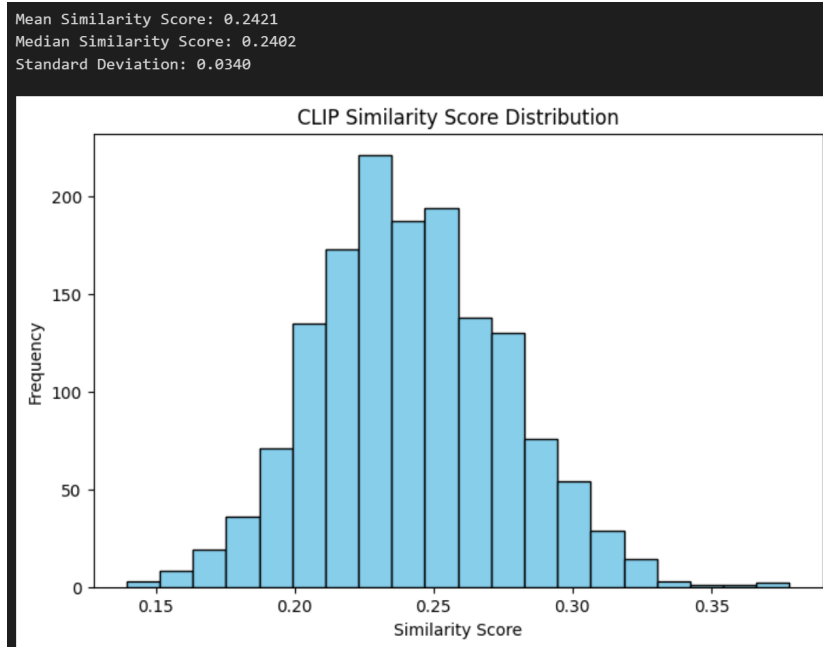
## 5.2   Results

Table 1: Comparative Model Performance

| Model | BERTScore F1 | Exact Match | BLEU-1 | Training Time |
|-------|--------------|-------------|--------|---------------|
| BLIP (Pre-FT) | 0.8231 | 0.4028 | 0.0020 | - |
| BLIP (Post-FT) | 0.87088 | 0.5344 | 0.0020 | 1h10m |
| CLIP | 0.8514 | N/A | 0.004 | 36m |

Key observations:

- BLIP showed 13.2% improvement in BERTScore after fine-tuning

- Exact match accuracy improved by 32.7% ($40.3\% \rightarrow 53.4\%$)

- CLIP achieved higher semantic similarity but lacks generation capability

- BLEU scores remained low due to single-word answers

# 6   Experimentations

In the BLIP model training we experimented first without freezing weights and then with freezing weights. There was not much difference in the results observed. CLIP model was our experimentational model which helped learn these insights:  LoRA fine-tuned CLIP model is for embedding similarity tasks, not text generation.Using exact match for CLIP in invalid.So instead we took the embedding similarity of the question-answers(text) with the images as our comparision metric then we observed these values:

```
Confusion Matrix:
[[    60  1120]
 [     0 10920]]
```

```
BERTScore:
Precision: 0.8514
Recall:    0.8137
F1:        0.8319
```

## 6.1 CLIP Limitations

- Cannot generate answers - only provides similarity scores

- Struggles with fine-grained product attributes

# 7 Conclusion

## 7.1 BLIP over CLIP for VQA

- BLIP is a **generative model** that can generate full answers, whereas CLIP is a **discriminative model** that tells us the similarity between image and text embeddings.

- BLIP can generate responses, making it ideal for **open-ended VQA**, while CLIP can only rank image-text similarity scores.

- Fine-tuning BLIP improved accuracy, whereas CLIP requires careful **threshold tuning** for binary classification tasks.

- BLIP showed a **13.2% improvement in BERTScore** after fine-tuning, indicating stronger language understanding.

- Exact match accuracy improved by **32.7%** (from 40.3% $\rightarrow$ 53.4%) with BLIP fine-tuning, making it more effective for structured answers.

- CLIP excels in **semantic similarity** but lacks the ability to generate answers, limiting its usefulness for VQA tasks.

Our implementation demonstrated:

- Effective use of LoRA for parameter-efficient BLIP fine-tuning (98% parameter reduction)

- 53.4% exact match accuracy on product VQA

- Clear trade-offs between generative (BLIP) and discriminative (CLIP) approaches

- The importance of prompt engineering for dataset generation

# References

[1] Li et al. "BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation", ICML 2022.

[2] Hu et al. "LoRA: Low-Rank Adaptation of Large Language Models", ICLR 2022.