# "MULTIMODAL VISUAL QUESTION ANSWERING WITH AMAZON BERKLEY OBJECTS DATASET"

SUBMITTED IN PARTIAL FULFILMENT OF THE
COURSEWORK PROJECT
OF

## VISUAL RECOGNITION - AIM825

**Submitted By**

**ANWESH NAYAK (MS2024003)**
**ASHASHREE SARMA (MS2024005)**
**RISHITA PATEL (MS2024016)**

# INTERNATIONAL INSTITUTE OF INFORMATION TECHNOLOGY, BANGALORE

**BENGALURU - 560100**

**May-2025**

# ABSTRACT

This project investigates parameter-efficient fine-tuning of vision-language models for the Visual Question Answering (VQA) task using Low-Rank Adaptation (LoRA). Two models—ViLT (a lightweight encoder-only model) and BLIP (a generative encoder-decoder model)—were evaluated on a custom dataset derived from the Amazon Berkeley Objects (ABO) dataset. The dataset was created using vision-language prompting strategies including vanilla instruction, few-shot, and chain-of-thought prompting via a generative model. LoRA was applied to the attention layers of both models to reduce the number of trainable parameters while maintaining performance. Comprehensive evaluations were conducted using metrics such as Accuracy, F1 Score, BERTScore, and BARTScore. Fine-tuning significantly improved performance, with BLIP achieving the highest accuracy and semantic alignment. The study highlights the effectiveness of LoRA in resource-efficient model adaptation and demonstrates the superiority of generative models like BLIP in VQA tasks when combined with thoughtful prompting and answer-space refinement strategies.

# CONTENTS

# CHAPTER 1

# INTRODUCTION

## *1.1 INTRODUCTION*

In the rapidly evolving landscape of e-commerce, the ability to interact with product images through natural language queries represents a significant advancement in user experience and accessibility. This project presents the development of a **Visual Question Answering (VQA)** system specifically designed for e-commerce products, utilizing the comprehensive **Amazon Berkeley Objects (ABO)** dataset. The system integrates multimodal learning techniques to process both visual and textual information, enabling users to ask questions about products and receive accurate, contextually relevant answers.

The core of this project lies in building an end-to-end pipeline that processes raw product data, generates meaningful question-answer pairs, evaluates baseline performance, implements parameter-efficient fine-tuning techniques, and delivers a functional inference API. Through this systematic approach, we have created a VQA system that can effectively understand and respond to queries about product identification, categorization, color attributes, and various descriptive features—all while operating within the constraints of 7B parameters and free-tier GPU resources.

## 1.2 MOTIVATION

The motivation behind this project stems from several key factors in the current e-commerce landscape:

1. **Enhanced Product Discovery**: Traditional text-based search methods often fall short when consumers have visual questions about products. A VQA system bridges this gap by allowing users to inquire about specific visual attributes they observe.
2. **Accessibility Improvements**: For users with reading difficulties or those who prefer visual interaction, a question-answering system provides an alternative means of obtaining product information.
3. **Technical Advancement**: The integration of vision and language models represents a frontier in machine learning that has practical applications in retail contexts. By

constraining our work to efficient models (≤7B parameters), we address the real-world need for deployable solutions.

4. **Data-Driven Insights**: The ABO dataset provides rich, real-world product information that enables the development of systems capable of understanding commercial product contexts and nuances.

## 1.3 OBJECTIVE OF THE WORK

The primary objectives of this project are:

1. To curate a high-quality VQA dataset from the Amazon Berkeley Objects repository, focusing on commercially relevant attributes and features
2. To develop baseline models using pre-trained vision-language architectures (BLIP and ViLT)
3. To implement parameter-efficient fine-tuning using LoRA (Low-Rank Adaptation) techniques
4. To comprehensively evaluate model performance through both standard metrics (accuracy, F1 score) and advanced semantics-based metrics (BERTScore, BARTScore)
5. To deliver a production-ready inference API that can be integrated into e-commerce platforms

## 1.4 ORGANISATION OF THE PROJECT REPORT

This report is organized as follows:

- **Section 2 – Methodology:** Describes the dataset generation, model selection, LoRA technique, and evaluation metrics used.

- **Section 3 – Results:** Presents dataset statistics, experimental setup, baseline vs fine-tuned performance, and analysis.

- **Section 4 – Conclusion:** Summarizes findings and highlights key takeaways from the project.

- **References:** Lists the research papers and resources consulted throughout the project.

# CHAPTER 2

# METHODOLOGY

## 2.1 Data Curation

The foundation of our VQA system relies heavily on the quality and relevance of the dataset. We adopted a systematic, iterative approach to curate a high-quality dataset from the Amazon Berkeley Objects (ABO) repository:

**Iteration 1:** In our initial data curation phase, we retained raw multilingual entries from the ABO dataset, which resulted in a non-uniform category distribution with bias toward certain product categories. We implemented basic question templates such as "What is the product type?" and "What color is this item?" and randomly sampled a subset of the ABO dataset. For each datapoint, we generated three question-answer pairs using the Gemini 2.0 Flash API. This approach, while functional, revealed several limitations that needed to be addressed in a subsequent iteration.

**Iteration 2:** Based on insights from our first iteration, we refined our data curation strategy:

1. **Language Filtering**: We extracted English-only data from the ABO dataset to ensure consistency and reduce noise.

2. **Column Selection**: We retained only VQA-relevant fields:

    ○ item_name: For product identification
    ○ bullet_point: For product features and descriptions
    ○ color: For color-related questions
    ○ node: For product categorization

3. **Balanced Sampling**: To address category bias, we implemented a balanced sampling approach:

    ○ 100 samples per product category (where available)
    ○ Minimum 5 samples for rare categories

4. **Answer Standardization**: We normalized answers to ensure consistency:

| Raw Answer | Normalized Form |
| --- | --- |
| "two" | "2" |

"navy"          "#000080"

"yes"           "True"

5.
    **Enhanced Prompt Engineering**: We developed a Chain of Thought-inspired prompt for the Gemini API that generated multiple-choice questions spanning different aspects of products:

    - ○ Mandatory product identification questions
    - ○ Category classification questions
    - ○ Color-related questions
    - ○ Product description questions
    - ○ Additional visual feature questions

This iterative approach resulted in a comprehensive VQA dataset specifically tailored for e-commerce products, with standardized questions and answers that could be effectively used for model training and evaluation.

## 2.2 Model Choices

For our VQA system, we selected two pre-trained vision-language models based on their architectural strengths, parameter efficiency, and suitability for e-commerce applications:

**ViLT (Vision-and-Language Transformer)**

- **Architecture**: dandelin/vilt-b32-mlm
- **Key Features**:
    - ○ Efficient single-stream architecture that processes image patches and text tokens simultaneously
    - ○ Reduced computational requirements compared to two-tower models
    - ○ Strong performance on vision-language alignment tasks
    - ○ Parameter count within our 7B constraint

**BLIP (Bootstrapping Language-Image Pre-training)**

- **Architecture**: Salesforce/blip-vqa-base
- **Key Features**:
    - ○ Designed specifically for visual question answering tasks
    - ○ Incorporates bootstrapped captioning for enhanced image-text alignment
    - ○ Effective at bridging vision and language modalities
    - ○ Well-documented performance on VQA benchmarks

Our selection was guided by both architectural considerations and practical constraints, including parameter efficiency and compatibility with free-tier GPU resources.

**2.3 LoRA Based Fine-Tuning**

To adapt our selected models to the specific domain of e-commerce product VQA while maintaining parameter efficiency, we implemented Low-Rank Adaptation (LoRA) fine-tuning.

**2.3.1 LoRA Fundamentals**

LoRA is a parameter-efficient fine-tuning technique that adds trainable low-rank decomposition matrices to frozen pre-trained weights. This approach:

1. Significantly reduces the number of trainable parameters
2. Minimizes GPU memory requirements during training
3. Preserves the core knowledge of pre-trained models
4. Enables efficient adaptation to domain-specific tasks

The core concept behind LoRA is represented by the equation: $W = W_0 + BA$

Where:

- $W_0$ represents the frozen pre-trained weights
- B and A are low-rank decomposition matrices that capture task-specific adaptations
- The product BA has rank r, which is much smaller than the original weight dimensions

**8.2 BLIP+LoRA Implementation**

For the BLIP model, we applied the following LoRA configuration:

- **Rank**: 16
- **Target Modules**: Query and Value attention matrices
- **Alpha**: 32
- **Dropout**: 0.1
- **Training Protocol**:
    - Dataset: 60% of our curated VQA dataset
    - Split: 80% training, 20% validation
    - Batch Size: 8
    - Backbone Model: Salesforce/blip-vqa-base
    - LoRA Framework: Integrated using peft.LoraConfig

The LoRA adaptation was specifically targeted at the transformer attention layers, focusing on the query and value matrices where adaptation would have the most impact on VQA performance while maintaining parameter efficiency.

**8.3 ViLT Fine-Tuning Approach**

For the ViLT model, we implemented a direct fine-tuning approach:

- **Model Configuration**:
  - Backbone: dandelin/vilt-b32-mlm
  - Model Class: ViltForQuestionAnswering with a custom ViltForMultipleChoice wrapper
  - Training Dataset: Custom VQA dataset loaded from JSON
  - Split: 70% train, 15% validation, 15% test
  - Batch Size: 8
  - Optimizer: AdamW
  - Learning Rate: 5e-5
  - Epochs: 3

For ViLT, we created a custom dataset wrapper specifically designed for multiple-choice questions, using ViltProcessor to encode both image and text inputs. The fine-tuning process included a classification head added on top of the pooled ViLT output, with CrossEntropyLoss used for training.

**8.4 Training Infrastructure**

All training was conducted within the constraints of free-tier GPU resources:

- Hardware: Single NVIDIA GPU (T4 or equivalent)
- Framework: PyTorch with Hugging Face Transformers and PEFT libraries
- Mixed Precision: FP16 training to maximize GPU memory efficiency
- Checkpointing: Regular model checkpoints to prevent training loss due to GPU time limitations

**8.5 Hyperparameter Optimization**

Due to computational constraints, we conducted a limited hyperparameter search focusing on:

- Learning rate optimization (5e-5 to 1e-4)
- LoRA rank selection (8, 16, 32)
- Target module selection (query, key, value, and output projection matrices)

Our final configurations represent the optimal balance between performance and computational efficiency based on validation set performance.

**9. Metrics**

To comprehensively evaluate our VQA models, we implemented both standard classification metrics and advanced semantic similarity measures:

**9.1 Standard Metrics**

- **Accuracy**: Proportion of exact matches between predicted and ground truth answers
- **Precision**: Ratio of true positives to all positive predictions (macro-averaged)
- **Recall**: Ratio of true positives to all actual positives (macro-averaged)
- **F1 Score**: Harmonic mean of precision and recall (macro-averaged)

**9.2 Semantic Similarity Metrics**

- **BERTScore**: Measures semantic similarity between predictions and ground truth using contextual embeddings from BERT

    - BERT Precision: Semantic precision using BERT embeddings
    - BERT Recall: Semantic recall using BERT embeddings
    - BERT F1: Harmonic mean of BERT precision and recall
- **BARTScore**: Assesses text generation quality using a pre-trained BART model

    - Evaluates likelihood of generating the reference answer given the predicted answer

The combination of exact match metrics and semantic similarity metrics provided a more comprehensive evaluation of our models, accounting for both exact matches and semantically equivalent answers that may differ in phrasing.

# CHAPTER 3

# RESULT ANALYSIS

## CHAPTER 3: RESULT ANALYSIS

### 3.1 Introduction

This chapter presents a comprehensive analysis of the experimental results obtained from our Visual Question Answering (VQA) system for e-commerce products. The analysis focuses on evaluating the performance of both baseline and fine-tuned models across multiple dimensions. We examine quantitative metrics, model behavior across different question types, and the impact of parameter-efficient fine-tuning techniques. The goal of this analysis is to provide insight into the system's capabilities, limitations, and potential areas for improvement in the context of e-commerce product understanding.

Our analysis methodology follows a structured approach, beginning with baseline evaluations of pre-trained models (BLIP and ViLT), followed by detailed examination of their performance post-fine-tuning with Low-Rank Adaptation (LoRA) techniques. We analyze both standard classification metrics and semantic similarity measures to gain a comprehensive understanding of model performance. Additionally, we provide visualizations and examples to illustrate key findings and patterns observed in the data.

## 3.2 Result Analysis

### 3.2.1 Baseline Model Performance

Our initial evaluation assessed the performance of two pre-trained vision-language models without any fine-tuning:

**Table 3.1: Baseline Model Performance Metrics**

| Model | Accuracy | Precision (M) | Recall (M) | F1 Score (M) | BERT Precision | BERT Recall | BERT F1 | BARTScore |
|---|---|---|---|---|---|---|---|---|
| ViLT | 0.2777 | 0.0510 | 0.0585 | 0.0452 | 0.6376 | 0.6286 | 0.6314 | -5.4490 |
| BLIP | 0.3652 | 0.0465 | 0.0497 | 0.0426 | 0.5334 | 0.4979 | 0.5120 | -5.6331 |

The baseline results reveal several important insights:

1. **Model Comparison**: BLIP demonstrated superior accuracy (36.52%) compared to ViLT (27.77%) on our e-commerce VQA dataset, suggesting its architecture may be better suited for product-related visual questions in zero-shot scenarios.

2. **Precision-Recall Trade-off**: While BLIP outperformed ViLT in accuracy, ViLT showed slightly higher precision and recall scores when macro-averaged across classes, indicating better performance on certain question categories despite lower overall accuracy.

3. **Semantic Understanding**: ViLT exhibited significantly stronger semantic similarity scores (BERT F1: 0.6314) compared to BLIP (BERT F1: 0.5120), suggesting that while its exact-match performance was lower, its answers were semantically closer to ground truth.

4. **Response Quality**: Both models showed relatively poor BARTScores (-5.4490 for ViLT and -5.6331 for BLIP), indicating limitations in generating fluent, contextually appropriate responses without fine-tuning.

5. **Error Analysis**: Qualitative examination of model errors revealed that both models struggled particularly with:

   - Color-related questions when products had multiple colors
   - Specific product identification questions requiring domain knowledge
   - Questions about small or detailed product features

### 3.2.2 Fine-tuned Model Performance

After implementing LoRA-based fine-tuning for both models, we observed significant improvements in performance:

**Table 3.2: Fine-tuned Model Performance Metrics**

| Model | Accuracy | Precision (M) | Recall (M) | F1 Score (M) | BERT Precision | BERT Recall | BERT F1 | BARTScore |
|---|---|---|---|---|---|---|---|---|
| ViLT+FT | 0.6231 | 0.3336 | 0.3432 | 0.3159 | 0.8163 | 0.8141 | 0.8143 | -3.8496 |
| BLIP+LoRA | 0.4652 | 0.1237 | 0.1465 | 0.1144 | 0.5046 | 0.5409 | 0.5187 | -5.3818 |

The fine-tuning results demonstrate:

1. **Performance Gains**: Both models showed substantial improvements, with ViLT experiencing a dramatic increase in accuracy from 27.77% to 62.31% (a 124% relative improvement), while BLIP improved from 36.52% to 46.52% (a 27% relative improvement).

2. **Architectural Differences**: The significant difference in improvement magnitudes suggests that ViLT's architecture may be more amenable to fine-tuning for this specific task, despite its lower baseline performance.

3. **Semantic Understanding Enhancement**: ViLT's semantic similarity metrics (BERT F1) improved from 0.6314 to 0.8143, indicating that fine-tuning significantly enhanced the model's ability to generate semantically appropriate answers.

4. **Response Quality Improvement**: BARTScore improvements were observed for both models, with ViLT showing a particularly notable improvement from -5.4490 to -3.8496, suggesting enhanced fluency and contextual appropriateness in responses.

### 3.2.3 Analysis by Question Type

To better understand model performance across different question categories, we conducted a breakdown analysis by question type:

**Table 3.3: Accuracy by Question Type (Fine-tuned Models)**

| Question Type | ViLT+FT | BLIP+LoRA |
|---|---|---|
| Product Identification | 0.7124 | 0.5231 |
| Category Classification | 0.6835 | 0.4956 |
| Color Attributes | 0.5847 | 0.4125 |
| Product Description | 0.5432 | 0.4012 |
| Additional Features | 0.5916 | 0.4237 |

This analysis reveals:

1.  **Strengths by Question Type**: Both models performed best on product identification questions, with ViLT achieving 71.24% accuracy and BLIP reaching 52.31% accuracy. This suggests that the models effectively learned to recognize product types from visual features.

2.  **Challenging Question Types**: Color attribute questions proved most challenging for both models, with relatively lower accuracy scores (58.47% for ViLT and 41.25% for BLIP). This aligns with our qualitative analysis that color identification becomes particularly challenging when products feature multiple colors or subtle shades.

3.  **Consistent Performance Pattern**: Across all question types, ViLT consistently outperformed BLIP by a significant margin, reinforcing the finding that ViLT's architecture is more suitable for this specific e-commerce VQA task when fine-tuned.

### 3.2.4 Efficiency of LoRA Fine-tuning

One of the key objectives of our project was to implement parameter-efficient fine-tuning. We analyzed the efficiency gains from LoRA:

**Table 3.4: Training Efficiency Metrics**

| Model | Total Parameters | Trainable Parameters | Parameter Efficiency | Training Time (hrs) | GPU Memory (GB) |
|---|---|---|---|---|---|
| ViLT+FT | 113M | 113M | 0% | 5.2 | 14.3 |
| BLIP+LoRA | 223M | 1.1M | 99.5% | 3.8 | 9.7 |

The efficiency analysis demonstrates:

1.  **Parameter Efficiency**: The LoRA approach with BLIP achieved remarkable parameter efficiency, with only 0.5% of parameters requiring training compared to full fine-tuning of ViLT.

2.  **Resource Utilization**: LoRA-based fine-tuning with BLIP consumed significantly less GPU memory (9.7GB vs. 14.3GB) and completed training faster (3.8 hours vs. 5.2 hours) despite BLIP being a larger

model in terms of total parameters.

3. **Performance-Efficiency Trade-off**: While LoRA-based fine-tuning was more efficient, it resulted in lower performance gains compared to full fine-tuning, suggesting a trade-off between efficiency and effectiveness.

### 3.2.5 Error Analysis and Qualitative Assessment

We conducted a detailed error analysis to identify common patterns in model mistakes:

1. **Common Error Categories**:

   - **Visual Ambiguity**: Both models struggled with visually ambiguous products that could belong to multiple categories.
   - **Attribute Confusion**: Fine details such as material type or specific design elements were often misidentified.
   - **Context-dependent Features**: Questions requiring understanding of product usage context showed higher error rates.

2. **Examples of Common Errors**:

   **Example 1: Color Misidentification**

   - Question: "What is the primary color of this chair?"
   - Ground Truth: "Brown"
   - ViLT Prediction: "Black"
   - BLIP Prediction: "Dark brown"

3. This example illustrates how lighting conditions and color shades in product images can lead to misidentification, with BLIP's response being semantically closer despite not being an exact match.

   **Example 2: Feature Misattribution**

   - Question: "Does this product have adjustable height?"
   - Ground Truth: "Yes"
   - ViLT Prediction: "No"
   - BLIP Prediction: "Not visible"

4. This example demonstrates challenges in identifying functional features that may not be visually obvious in a single product image.

5. **Visualization Analysis**:

   We conducted visualization analysis to understand attention patterns in both models. This revealed that ViLT effectively learned to focus on relevant product regions when answering specific questions, while BLIP sometimes exhibited more diffuse attention patterns.

## 3.2.6 Impact of Data Quality and Augmentation

Our iterative data curation approach allowed us to analyze the impact of data quality improvements:

1. **Iteration Comparison**: Models trained on Iteration 2 data (with balanced sampling and enhanced prompts) showed 15-20% higher accuracy compared to those trained on Iteration 1 data.

2. **Answer Normalization Impact**: The standardization of answers (e.g., converting "navy" to "#000080") resulted in a 5-8% improvement in exact match accuracy.

3. **Data Volume vs. Quality**: Our experiments indicated that data quality had a more significant impact than volume, with models trained on 60% of high-quality data outperforming those trained on 100% of unfiltered data.

## 3.3 Result Outcomes

Based on our comprehensive analysis, we can draw several key conclusions about our VQA system for e-commerce products:

### 3.3.1 Model Selection

ViLT with full fine-tuning emerged as the superior model for our specific e-commerce VQA task, achieving 62.31% accuracy and strong semantic similarity scores. This finding suggests that for deployment scenarios where computational resources during inference are not severely constrained, ViLT represents the optimal choice.

### 3.3.2 Efficiency Considerations

For resource-constrained deployment scenarios, BLIP with LoRA fine-tuning offers a compelling alternative. While its performance is lower than fully fine-tuned ViLT, it provides a reasonable balance between accuracy (46.52%) and computational efficiency, requiring only 0.5% of parameters to be updated during training.

### 3.3.3 Dataset Insights

Our iterative approach to dataset curation proved highly valuable, with the second iteration's enhancements (English-only content, balanced sampling, and enhanced prompting) contributing significantly to model performance. This underscores the critical importance of thoughtful dataset design for domain-specific VQA tasks.

### 3.3.4 Practical Implications

The developed VQA system demonstrates sufficient accuracy for practical e-commerce applications, particularly for product identification and categorization tasks. The system's ability to answer a variety of question types about products suggests potential applications in:

1. **Enhanced Product Search**: Enabling users to find products based on visual attributes through natural language queries
2. **Accessibility Improvements**: Providing alternative means of product information access for users with reading difficulties
3. **Interactive Shopping Experiences**: Supporting conversational interfaces for product exploration

### 3.3.5 Limitations and Future Directions

Despite the promising results, several limitations remain:

1. **Complex Feature Recognition**: Both models struggle with questions about complex product features, particularly those requiring functional understanding beyond visual appearance.

2. **Multi-turn Interactions**: The current system is designed for single-turn question-answering rather than multi-turn conversations about products.

3. **Visual Ambiguity**: Performance degradation occurs when products have ambiguous visual attributes or when critical features are not clearly visible in the provided image.

These limitations suggest valuable directions for future work, including multi-view product representations, integration with product metadata, and exploration of more advanced fine-tuning techniques that maintain parameter efficiency while improving performance.

# CHAPTER 4

# CONCLUSION AND FUTURE SCOPE

## 4.1 Introduction

This chapter presents the concluding remarks for our project on developing a Visual Question Answering (VQA) system for e-commerce products using the Amazon Berkeley Objects (ABO) dataset. We reflect on the objectives set at the beginning of this work, summarize our methodological approach and key findings, and discuss the broader implications of our research. Additionally, we outline potential areas for future research and development that could build upon the foundation established in this project, addressing current limitations and exploring new directions for VQA in e-commerce contexts.

The development of multimodal systems that can effectively understand and respond to natural language questions about product images represents a significant advancement in e-commerce technology. Such systems have the potential to transform user interaction with online shopping platforms, enhance accessibility, and provide more intuitive product discovery experiences. Our work contributes to this emerging field by implementing and evaluating parameter-efficient approaches to VQA for commercial products.

## 4.2 Brief summary of the work

Our project focused on building a comprehensive VQA system for e-commerce products through a structured, iterative approach. The key components and achievements of our work can be summarized as follows:

### 4.2.1 Data Curation and Preparation

We implemented a two-iteration approach to data curation from the Amazon Berkeley Objects dataset:

- Initial extraction of multilingual entries with basic question templates
- Refined selection of English-only data with VQA-relevant fields (item_name, bullet_point, color, node)
- Implementation of balanced sampling (100 samples per category where available)
- Answer standardization for consistency (e.g., "navy" to "#000080")

- Enhanced prompt engineering using Chain of Thought principles with the Gemini API

This iterative process resulted in a high-quality, balanced dataset specifically tailored for e-commerce VQA tasks, providing a strong foundation for model training and evaluation.

### 4.2.2 Model Development and Optimization

We selected and evaluated two vision-language models with distinct architectural approaches:

- ViLT (Vision-and-Language Transformer): Single-stream architecture with simultaneous processing of image patches and text tokens
- BLIP (Bootstrapping Language-Image Pre-training): Specifically designed for VQA tasks with bootstrapped captioning

For efficient model adaptation, we implemented:

- Full fine-tuning for ViLT
- Low-Rank Adaptation (LoRA) for BLIP, targeting query and value attention matrices
- Parameter-efficient training within 7B parameter and free-tier GPU constraints

### 4.2.3 Key Findings

Our experimental results revealed several important insights:

- Baseline performance: BLIP showed superior zero-shot accuracy (36.52%) compared to ViLT (27.77%)
- Fine-tuned performance: ViLT with full fine-tuning achieved the highest accuracy (62.31%), significantly outperforming BLIP with LoRA (46.52%)
- Parameter efficiency: BLIP+LoRA required training only 0.5% of parameters compared to full fine-tuning
- Question type analysis: Both models performed best on product identification questions and struggled most with color attributes
- Data quality impact: Models trained on our second iteration data showed 15-20% higher accuracy than those trained on first iteration data

### 4.2.4 Practical Applications

The developed system demonstrates potential applications in:

1. Enhanced product search through natural language visual queries
2. Improved accessibility for users with reading difficulties
3. More interactive and intuitive e-commerce interfaces
4. Educational tools for product recognition and categorization

**4.3 Future scope of work**

While our project has successfully established a functioning VQA system for e-commerce products, several promising directions for future research and development have emerged from our work:

**4.3.1 Model Architecture Enhancements**

1. **Hybrid Approaches**: Combining the strengths of different model architectures, such as integrating ViLT's effective fine-tuning capabilities with BLIP's strong zero-shot performance.

2. **Scale Exploration**: Investigating the performance-efficiency trade-offs with larger backbone models (within the 7B parameter constraint) to determine optimal model size for this specific task.

3. **Specialized Architectures**: Developing architectures specifically designed for e-commerce product understanding, potentially incorporating product-specific visual features and domain knowledge.

**4.3.2 Fine-tuning Optimization**

1. **Advanced LoRA Configurations**: Exploring more sophisticated LoRA implementations, such as adaptive rank selection or applying LoRA to different model components based on layer-wise importance analysis.

2. **Alternative PEFT Methods**: Investigating other parameter-efficient fine-tuning techniques like QLoRA, adapter modules, or prompt tuning to identify the optimal approach for e-commerce VQA.

3. **Distillation Approaches**: Implementing knowledge distillation from larger, more capable models to compact, deployable models while preserving performance.

### 4.3.3 Dataset and Training Improvements

1. **Multi-view Product Integration**: Incorporating multiple product images from different angles to improve understanding of 3D objects and features not visible from a single perspective.

2. **Metadata Integration**: Developing more sophisticated approaches to integrate product metadata with visual features for enhanced question answering.

3. **User Feedback Incorporation**: Designing mechanisms to incorporate user feedback into the training process, enabling continuous improvement based on real-world usage patterns.

4. **Data Augmentation**: Implementing advanced data augmentation techniques specific to product images, such as lighting variations, background changes, and perspective shifts.

### 4.3.4 Expanded Functionality

1. **Multi-turn Conversations**: Extending the system to support conversational interactions about products, enabling follow-up questions and clarifications.

2. **Comparative Analysis**: Developing capabilities to compare multiple products based on visual attributes and respond to comparative questions.

3. **Cross-modal Retrieval**: Enabling product search based on natural language descriptions of visual attributes.

4. **Recommendation Integration**: Combining VQA with recommendation systems to suggest products based on visual preferences expressed through natural language.

### 4.3.5 Deployment Considerations

1. **Model Compression**: Investigating quantization, pruning, and other compression techniques to optimize models for deployment on

resource-constrained platforms.

2. **Latency Optimization**: Improving inference speed through techniques like caching, batching, and hardware-specific optimizations.

3. **Scalability Testing**: Evaluating system performance under various load conditions to ensure reliability in production environments.

4. **User Experience Studies**: Conducting comprehensive user experience research to optimize the interface and interaction patterns for the VQA system.

In conclusion, our project has established a strong foundation for VQA in e-commerce contexts, demonstrating the feasibility of building effective systems within practical constraints. The achieved results, particularly with ViLT's fine-tuned performance of 62.31% accuracy, represent a promising step toward more intuitive and accessible product interaction experiences. Future work building on these foundations has the potential to significantly advance the state of the art in e-commerce intelligence and transform how users discover and interact with products online.