

Assignment 4 Report

Student Name: Rosie Wang, 1806394

Report: Assignment 4

Examples of Generated Text Before and After Fine-tuning

Below are 5 representative examples comparing captions from the pre-trained model and the LoRA fine-tuned model:

Example 1: Dogs with Toy

Ground Truth: One white and tan dog is chasing another white and brown dog who has a toy in its mouth.

Pre-trained: two dogs chasing each other on a beach

Fine-tuned: dogs are being chased by a white dog carrying a yellow toy. white and brown sand dune

Ground Truth: One white and tan dog is chasing another white and brown dog who has a toy in its mouth .

Pre-trained: two dogs chasing each other on a beach

Fine-tuned: dogs are being chased by a white dog carrying a yellow toy . white and brown sand dune



Example 2: Dogs Playing

Ground Truth: two tan dogs play in the grass

Pre-trained: a dog that is standing in the grass

Fine-tuned: tan dogs play in the grass. black and white are blurred. white is the background. brown

Ground Truth: two tan dogs play in the grass

Pre-trained: a dog that is standing in the grass

Fine-tuned: tan dogs play in the grass . black and white are blurred . white is the background . brown



Example 3: Frisbee Catch

Ground Truth: A dog jumps to catch an orange Frisbee.

Pre-trained: a dog jumping in the air to catch a frisbee

Fine-tuned: black dog is catching a red Frisbee in its mouth. blue and white, a brown

Ground Truth: A dog jumps to catch an orange Frisbee .

Pre-trained: a dog jumping in the air to catch a frisbee

Fine-tuned: black dog is catching a red Frisbee in its mouth . blue and white , a brown



Example 4: Playing with Toy

Ground Truth: An older dog and a younger one playing with a toy.

Pre-trained: two dogs playing with each other in the snow

Fine-tuned: tan dogs are playing tug of war with a brown ball. brown dog is trying to get it

Ground Truth: An older dog and a younger one playing with a toy .

Pre-trained: two dogs playing with each other in the snow

Fine-tuned: tan dogs are playing tug of war with a brown ball . brown dog is trying to get it



Example 5: Running on Shore

Ground Truth: A white and brown dog runs along the shoreline.

Pre-trained: a dog running across a sandy beach

Fine-tuned: wet white dog running along side a rocky shoreline. reflections in the distance.

Ground Truth: A white and brown dog runs along the shoreline .

Pre-trained: a dog running across a sandy beach

Fine-tuned: wet white dog running along side a rocky shoreline . reflections in the distance .



Key Observation: The fine-tuned model produces more detailed and specific descriptions that better match the ground truth style, using dataset-specific vocabulary like "tan dogs," "shoreline," and action-specific phrases like "tug of war."

Short Reflection(difference or theory from gemini)

Summary:

- **Dataset:** Flickr30k (50 images)
- **LoRA Config:** r=32, α =64, targets=[c_attn, c_proj, c_fc]
- **Training:** 8 epochs, ~10-15 minutes, learning rate 3e-4
- **Results:** Clear visible differences in 5/5 examples, better alignment with ground truth style and

Evaluating CLIP zero-shot on 10000 images...

CLIP Evaluation: 100%

40/40 [00:23<00:00, 2.24it/s]

=====

CLIP Zero-shot Test Accuracy: 88.85%

=====

Model Comparison on CIFAR-10:

CLIP (Zero-shot):	88.85%
ViT (Trained):	70.63%
CNN (Trained, A3):	10.35%

Observations:

- CLIP achieves 88.85% accuracy without any training on CIFAR-10
- This demonstrates the power of vision-language pre-training
- CLIP can generalize to new visual concepts through text descriptions
- Trained models (ViT/CNN) are optimized specifically for CIFAR-10

Why LoRA is Useful

LoRA (Low-Rank Adaptation) proves exceptionally valuable for student-scale experiments by enabling fine-tuning of large models with minimal computational resources. In this assignment, we successfully adapted a vision-language model using only 1.6% of its parameters (~2M out of 124M), completing training in 10-15 minutes on a free Google Colab GPU. This efficiency is achieved through low-rank matrix decomposition, where trainable matrices are injected into frozen pre-trained weights, allowing domain adaptation without requiring expensive hardware or extensive training time. The resulting LoRA adapter is tiny (just 8MB), making it easy to share and version control, while the frozen base model prevents catastrophic forgetting of general knowledge. This combination of parameter efficiency, fast training, and preserved base capabilities makes LoRA ideal for educational settings where students can experiment with state-of-the-art fine-tuning techniques despite limited computational budgets (from Gemini).

How Fine-tuning Changes Model Behavior

Fine-tuning with LoRA fundamentally transforms the model's generation behavior by adapting it to specific patterns in the target dataset. As our results demonstrate, the pre-trained model generates generic descriptions like "a dog running across a sandy beach," while the fine-tuned version produces more elaborate captions such as "wet white dog running along side a rocky shoreline. reflections in the distance," showing clear adaptation to Flickr30k's descriptive style. This transformation occurs through several mechanisms: the model learns dataset-specific vocabulary (e.g., "tan dogs," "shoreline"), develops more detailed attribute descriptions (colors, textures, actions), and adjusts its generation length and complexity to match the training data. By targeting both attention layers (`c_attn`, `c_proj`) and MLP layers (`c_fc`) with our LoRA configuration (rank $r=32$, alpha $\alpha=64$), we enable comprehensive adaptation of both attention patterns and feature transformations. Our aggressive training settings (learning rate $3e-4$, 8 epochs) ensure visible changes, as evidenced by the noticeably different outputs where the fine-tuned model attempts more specific color descriptions, detailed actions like "tug of war," and scene elements like "reflections in the distance," effectively creating a specialized version that better serves our image captioning domain.