

# Enhancing Visual Captioning: Comparing LoRA Fine-Tuning and Few-Shot Prompting on the PixelProse Dataset

Kyle Merritt  
University of Central Florida  
CAP 5415: Computer Vision

Fall 2025

## Abstract

Modern vision-language models can generate fluent image captions but often default to short, generic descriptions that underspecify relationships, context, and fine-grained details. This project investigates two strategies for improving detailed captioning: (1) parameter-efficient fine-tuning using Low-Rank Adaptation (LoRA) on a subset of the PixelProse dataset, and (2) multimodal few-shot prompting that conditions the model on exemplar image-caption pairs at inference time. Using the open-source InternVL 3.5 2B model as a base, I compare six configurations that vary both model state (pretrained versus LoRA-fine-tuned) and prompting style (neutral, detailed, and detailed with multimodal few-shot exemplars). Quantitative evaluation on 500 held-out PixelProse images uses BLEU, METEOR, CIDEr, BERTScore, and caption length. Results show that multimodal few-shot prompting substantially improves similarity to PixelProse captions for the pretrained model, while the limited LoRA fine-tuning run used in this study yields only modest additional gains. Combining LoRA with few-shot prompting degrades performance and produces repetitive, off-target captions, which I attribute to a mismatch between the single-image training prompt and the multi-image few-shot evaluation prompt. I also measure the additional VRAM cost of multimodal few-shot prompting and discuss the trade-off between inference-time memory, training cost, and caption quality.

## 1 Introduction

Dense, descriptive image captioning is important for applications such as assistive technologies, visual storytelling, and downstream reasoning in multimodal agents. In practice, however, many off-the-shelf vision-language models (VLMs) still produce short, generic captions that under-describe complex scenes (for example, “a woman wearing a jacket” instead of specifying pose, attributes, and context). Improving caption richness can be approached either by modifying the model itself through fine-tuning or by reformulating the prompts used at inference time.

In this project I investigate these two strategies in the context of InternVL 3.5 2B, a recent open-source multimodal model, on the PixelProse dataset of dense captions. The first strategy is using Low-Rank Adaptation (LoRA), a parameter-efficient fine-tuning method that keeps all original Transformer weights frozen and injects small trainable low-rank matrices into existing attention and MLP projection layers, allowing the model to specialize to PixelProse-style captions without modifying or adding full layers. The second strategy is multimodal few-shot prompting: I prepend exemplar image-caption pairs to the model’s input so that it can condition on several high-quality examples before describing a new query image.

The original project proposal envisioned training on millions of PixelProse examples with multi-GPU data-parallel fine-tuning and then running a large ablation over LoRA ranks, frozen-layer depths, and several few-shot prompting schemes. In practice, I was constrained to single-GPU training and a much smaller subset of the dataset, and I evaluated each configuration on 500 held-out images instead of the larger planned study. Despite these limitations, the experiment still answers three central questions:

- How effective is multimodal few-shot prompting for improving descriptive caption quality in this domain?
- How much additional benefit does LoRA fine-tuning provide beyond prompting when training is performed on a relatively small subset?
- Do fine-tuning and few-shot prompting combine synergistically, or can mismatched training and inference conditions actually hurt performance?

## 2 Related Work

### 2.1 Dense Captioning and PixelProse

Traditional captioning datasets such as MSCOCO emphasize short alt-text style captions that summarize only the most salient objects in an image. In contrast, PixelProse [3] introduces more than 16 million image–caption pairs featuring long, dense, paragraph-level descriptions intended to capture fine-grained attributes, relationships, and scene structure. Each record also includes metadata such as aesthetic scores and watermark probabilities, enabling quality-based filtering and analysis. Because PixelProse captions are significantly longer and more detailed than standard alt-text, the dataset provides a natural benchmark for studying models that aim to generate richly descriptive prose.

### 2.2 InternVL 3.5

InternVL 3.5 [4] is an open-source family of multimodal models that integrate a vision encoder with a large language model. The 2B-parameter variant used in this project is small enough for single- or dual-GPU setups while still supporting multi-image reasoning, conversation-style interactions, and flexible input formats. InternVL employs dynamic image tiling to handle high-resolution or non-square images by splitting them into patches, encoding each patch through its vision backbone, and injecting learned image tokens into the language model. This architecture enables the model to process varying image sizes without retraining.

### 2.3 LoRA Fine-Tuning

Low-Rank Adaptation (LoRA) [2] is a parameter-efficient fine-tuning method that freezes the original model weights and injects trainable low-rank matrices into targeted linear layers. Instead of updating full projection or feed-forward matrices, LoRA learns small rank- $r$  updates that are added to the existing weights at inference time. This approach drastically reduces memory footprint and training compute, making it well suited for adapting large models on commodity hardware. In this project, LoRA is applied to the transformer blocks of the language model component of InternVL by targeting key attention and MLP projection layers; only the LoRA parameters are updated while the base model remains unchanged.

### 2.4 Few-Shot Prompting for Vision–Language Models

Few-shot prompting has long been used in large language models by providing example input–output pairs directly within the prompt. This idea extends naturally to multimodal models by including paired image–caption exemplars. A landmark demonstration is Flamingo [1], which showed that large pretrained vision–language models can achieve strong few-shot performance on captioning, VQA, and multimodal reasoning tasks simply by concatenating exemplar image–caption pairs with the query input, without any parameter updates. Such prompting enables rapid style conditioning—such as inducing longer, more descriptive captioning behavior—but increases inference-time VRAM usage because all exemplar images must be embedded and stored in the attention context.

In this project, we evaluate whether multimodal few-shot prompting can effectively guide InternVL toward producing PixelProse-style long-form captions and compare its benefits and costs to those of LoRA fine-tuning.

## 3 Method

### 3.1 Dataset and Subset Construction

The full PixelProse dataset is too large to train on exhaustively within a single semester project, so I constructed several smaller subsets using a dedicated script (`scripts/build_pixelprose_subset.py`). The script downloads images and metadata from the Hugging Face PixelProse dataset, shuffles with a fixed seed, filters out failed downloads, and writes each subset to a directory of the form:

`data/pixelprose_subset{index}/`

Each subset contains:

- An `images/` directory with JPEG files named `000000.jpg`, `000001.jpg`, ...
- A `metadata.jsonl` file with one JSON record per image, including fields such as `id`, `url`, `image_file`, `caption`, `vlm_model`, aesthetic score, and watermark score.

For the main experiment I used **subset 2** (`data/pixelprose_subset2`), which contains 2,000+ images. I reserved all examples with `id < 1000` for evaluation and used `id ≥ 1000` for LoRA training and validation. The LoRA training script further splits this usable pool into a training set and a validation set using a fixed, deterministic partition.

In practice, due to runtime and memory constraints, I trained on a much smaller number of images than originally proposed. The final LoRA run used approximately 20,000 captioned examples sampled from the training split through multiple epochs, and I evaluated each configuration on 500 held-out images from the evaluation pool (except where otherwise noted).

### 3.2 Base Model and Image Preprocessing

All experiments use the `OpenGVLab/InternVL3.5-2B-Instruct` checkpoint. Images are preprocessed using the dynamic tiling pipeline implemented in the load image function within `src/test_internvl_caption.py`, which follows the official InternVL preprocessing strategy. The steps are:

1. Load the image from disk using PIL and convert it to RGB.
2. Use `dynamic_preprocess` to select an aspect-ratio-appropriate grid of tiles, producing up to 12 base patches; if more than one tile is used, a  $448 \times 448$  thumbnail patch is appended, giving a maximum of 13 patches.
3. Resize each patch to  $448 \times 448$  using bicubic interpolation and normalize using ImageNet mean and standard deviation.
4. Stack the patches into a tensor of shape  $[N_{\text{patch}}, 3, 448, 448]$ .

This preprocessing function is shared across both fine-tuning and evaluation, ensuring that the visual input distribution is consistent at training and test time.

### 3.3 LoRA Configuration and Training Pipeline

LoRA fine-tuning is implemented in `src/train_lora.py`. The script loads the full `OpenGVLab/InternVL3.5-2B-Instruct` checkpoint and applies a LoRA adapter to the `language_model` component inside the InternVL architecture. The configuration is:

- Rank  $r = 32$ , LoRA scaling  $\alpha = 64$ , dropout = 0.05.
- Target modules: the query, key, value, and output projections (`q_proj`, `k_proj`, `v_proj`, `o_proj`), as well as the gate, up, and down projections in the feed-forward layers (`gate_proj`, `up_proj`, `down_proj`).
- Task type: causal language modeling, with cross-entropy loss over the full token sequence (image tokens, prompt, and caption).

The training dataset class `PixelProseLoraDataset` constructs each example from a single image and its PixelProse caption. It uses the same detailed prompt as the PT/FT evaluation setting (`prompts/detailed.txt`):

*“Provide a detailed, comprehensive caption describing all key objects, attributes, actions, and context.”*

For each training sample, the multimodal input fed to the model is built in `build_mm_inputs_for_batch`: a block of special image tokens `<img> ...<IMG.CONTEXT> ...</img>` is generated, where the `<IMG.CONTEXT>` token is repeated in proportion to the number of image patches for that example. This image-token block is followed by the detailed prompt and finally the ground-truth caption. No few-shot exemplars are included during training; every example consists of a single query image and its caption.

The collate function concatenates all image patches in a batch into a single tensor and tracks the number of patches per sample as `num_patches_list`. This allows the training loop to (i) move all patches onto the GPU as one tensor and (ii) compute the correct number of visual tokens for each query when constructing the text sequence.

All experiments were run on a single RTX 3090 GPU. The training script is configured for single-GPU operation (`device_map = None`) and uses both a limited number of epochs and a hard cap on the total number of optimization steps (`max_steps`) to keep runs within the available compute budget. As a result, the LoRA adapter in this project should be viewed as a relatively small fine-tuning run rather than a fully converged model. I attempted to enable multi-GPU LoRA fine-tuning using HuggingFace Accelerate, but InternVL’s multimodal architecture cannot be automatically sharded across devices without modifying the underlying model; because the model does not implement HF’s parallelization hooks and contains custom cross-modal injection layers, multi-GPU training was not feasible within the project window.

By default, the adapter weights are saved under `checkpoints/internvl3.5_2b.lora.pixelprose/subset2_r32_a64/`. In my experiments, I also produced a variant checkpoint with an explicit training-size suffix (`subset2_r32_a64_train_size_20000`). The evaluation script `src/eval_captions.py` loads the LoRA adapter path (either the default or the suffixed directory) when running the fine-tuned (FT) configurations.

### 3.4 Prompting Configurations

Prompting is controlled by three instruction templates stored in the `prompts/` directory:

- **Neutral (N):** “Provide a factual caption for this image.”
- **Detailed (D):** “Provide a detailed, comprehensive caption describing all key objects, attributes, actions, and context.”
- **Detailed Few-Shot (D\_FS):** “Provide a detailed, comprehensive caption describing all key objects, attributes, actions, and context in the query image. Use the exemplar image–caption pairs above as a guide for the level of detail and style, but focus entirely on the query image.”

The main evaluation script `src/eval_captions.py` supports both pretrained (PT) and fine-tuned (FT) models and can evaluate multiple prompting methods on the same set of images. For this project, I focused on the following six configurations:

Model state	Prompting method
PT	N (neutral)
PT	D (detailed)
PT	D_FS (detailed + few-shot)
FT	N (neutral)
FT	D (detailed)
FT	D_FS (detailed + few-shot)

Few-shot prompting uses multimodal exemplars drawn from the same subset. For each run, I select a small, fixed set of exemplar records from the metadata file and log them to

`logs/eval_subset{idx}-{PT/FT}_exemplars.jsonl`. To avoid test leakage, I ensure that these exemplar IDs are excluded from the evaluation pool and verify this with a dedicated script (`scripts/check_exemplar_leakage.py`).

In the multimodal few-shot setup, I prepend  $K$  exemplar image–caption pairs (typically  $K = 3$ ) before the query image. At evaluation time, the text prompt looks conceptually like:

```
Exemplar-1: <image>
Caption-1:  [caption text]
...
Query: <image>
[detailed few-shot instruction]
```

This design allows the model to see multiple PixelProse-style captions before generating a description for the query image.

The three exemplar image caption pairs that were used in the few shot prompting configuration are shown below.



**Exemplar 1 (outfit image, ID 0):**

This image displays two sets of outfits, each containing a black T-shirt, a blue blazer, and a pair of blue jeans. The T-shirts have a white logo on them, reading: "S.T.A.R. Labs." The blazers are both a deep royal blue color. The left blue blazer has a white lightning bolt emblem on the left sleeve. The right blazer has a gold emblem of an atom on the left sleeve. Both pairs of jeans are light blue and slightly distressed. The left outfit is paired with a pair of red high-top Converse sneakers, and a red purse with black edges and handles. The right outfit is paired with a pair of black high-top Converse sneakers and a silver necklace with a lightning bolt pendant. The background of the image is a white square.

**Metadata:**

Aesthetic score: 5.30    Watermark score: 0.83

**Role in prompts:**

This image and caption pair is used as the first exemplar in the multimodal few-shot prompt (PT\_D\_FS and FT\_D\_FS). It provides a long, style-rich description of a fashion-oriented image, demonstrating the level of detail and structure that the model is encouraged to imitate.

Figure 1: Exemplar 1 used in the multimodal few-shot prompt: a detailed fashion style image and caption pair.



**Exemplar 2 (crest image, ID 1):**

This image displays:

A crest with a shield divided into four sections. A lion, a snake, a badger, and a raven are in each section, respectively. The letters “H” and “P” are in the middle of the crest. A banner with the words “Draco Dormiens Nunquam Titillandus” is at the bottom. The crest has a gold background and a gold frame with flourishes.

**Metadata:**

Aesthetic score: 5.61    Watermark score: 0.90

**Role in prompts:**

This is the second exemplar in the multimodal few-shot prompt. It shows that PixelProse style captions can include emblematic or symbolic content with dense lexical detail, including text, heraldry elements, and layout information.

Figure 2: Exemplar 2 used in the multimodal few-shot prompt: a detailed crest image and caption pair.



**Exemplar 3 (cartoon quote image, ID 2):**

This image displays a pink background with black text. The text reads: “I just want to lose weight while staying in bed, watching TV, and eating Girl Scout cookies. Is that really too much to ask?” The text is in a white speech bubble on the left side of the image. On the right side of the image is a cartoon drawing of a person sitting on a couch with polka dots. The person is wearing pajamas and has their head in their hands. There is a box of Girl Scout cookies on the floor in front of the couch. The image is drawn in a simple cartoon style.

**Metadata:**

Aesthetic score: 5.18    Watermark score: 0.85

**Role in prompts and failure analysis:**

This is the third exemplar in the multimodal few-shot prompt. Its caption is later mistakenly reproduced almost verbatim by the FT + D\_FS configuration when captioning a different query image (Figure 9), illustrating a clear case where the fine-tuned model copies an exemplar instead of describing the query.

Figure 3: Exemplar 3 used in the multimodal few-shot prompt: a long, stylized caption for a cartoon quote image. This exemplar is directly involved in the exemplar confusion failure shown in Figure 9.

Table 1: Quantitative results on 500 held-out PixelProse images from subset 2 (451 for FT + D\_FS). Average lengths are in tokens.

Model	Prompt	#	BLEU	METEOR	CIDEr	BERTScore
PT	N	500	3.87	18.66	0.09	40.25
PT	D	500	8.80	28.97	0.18	43.20
PT	D_FS	500	<b>10.08</b>	<b>33.05</b>	<b>0.19</b>	<b>47.05</b>
FT	N	500	8.30	32.12	0.04	42.67
FT	D	500	8.69	33.41	0.04	43.86
FT	D_FS	451	7.06	30.60	0.01	38.58

### 3.5 Evaluation Protocol and Metrics

For each configuration, the evaluation script runs InternVL on up to 500 held-out images from `pixelprose_subset2`, producing logs of the form:

```
logs/eval_subset2_{PT/FT}_{N,D,D_FS}.jsonl
```

Each record includes the image ID, file name, ground-truth caption, prompt type, model tag (PT or FT), and the generated prediction.

Automatic metrics are computed by `scripts/compute_metrics.py`, which aggregates:

- BLEU (via `sacrebleu`),
- METEOR (via `NLTK`),
- CIDEr (via a Hugging Face `evaluate` implementation),
- BERTScore (F1, rescaled to 0–100),
- Average ground-truth and predicted caption lengths (in tokens).

The main summary for subset 2 and the PT/FT configurations is stored in `metrics/metrics_subset2_PT_FT.csv`. For most configurations I obtained results on all 500 examples. The FT + D\_FS configuration was ended early because multiple repetitions and incorrect descriptions were observed, resulting in 451 evaluated samples.

In addition to quantitative metrics, I manually reviewed generated captions to identify typical successes and failure modes. To make these qualitative differences visible, I prepared galleries pairing each image with its PixelProse caption and the generated captions under several configurations. Section 4.3 describes how these galleries are organized in the report.

## 4 Experiments and Results

### 4.1 Quantitative Results

Table 1 summarizes the main automatic metrics for the six configurations on subset 2. All metrics are reported on a 0–100 scale for easier comparison.

Several trends are immediately apparent:

- For the pretrained model, simply switching from the neutral prompt (PT + N) to the detailed prompt (PT + D) yields a large improvement across all metrics. Average predicted caption length increases from roughly 40 tokens to nearly 100 tokens, and BLEU, METEOR, CIDEr, and BERTScore all improve.
- Adding multimodal few-shot exemplars on top of the detailed prompt (PT + D\_FS) produces the best overall scores for the pretrained model. METEOR and BERTScore in particular increase noticeably, suggesting that the model’s style and content become more aligned with PixelProse captions when it sees exemplar image–caption pairs.

- LoRA fine-tuning with the neutral prompt (FT + N) dramatically increases average caption length (often exceeding 190 tokens) and provides higher METEOR than PT + N. However, CIDEr remains low, indicating that the increased verbosity does not always translate into better n-gram agreement with the ground truth.
- The fine-tuned model with the detailed prompt (FT + D) performs similarly to PT + D in terms of BLEU and METEOR, with slightly higher BERTScore but still modest gains compared to the impact of multimodal few-shot prompting.
- Surprisingly, combining LoRA fine-tuning with detailed few-shot prompting (FT + D\_FS) *hurts* performance. All metrics drop relative to PT + D\_FS, and BERTScore decreases substantially. Qualitative inspection reveals that this configuration often produces repetitive or off-topic captions.

Overall, on this training budget, multimodal few-shot prompting provides the clearest benefit over the base model, whereas the small-scale LoRA fine-tuning run yields only modest additional improvements and interacts poorly with few-shot prompting.

## 4.2 VRAM Usage and Computational Cost

One of the original goals of this project was to compare not only caption quality but also computational cost. Although I did not perform a full runtime and energy audit, I did track GPU memory usage during key configurations.

When using multimodal few-shot prompting with three exemplar images plus one query image, InternVL must process all patches from four images in a single forward pass. On my hardware, this configuration used approximately 50 GB of VRAM across the 4x RTX 3090 system. This is significantly higher than the memory required for single-image captioning or for running the fine-tuned model without exemplars.

In contrast, LoRA fine-tuning itself can be performed with more modest per-batch memory (once the model is loaded), and the resulting fine-tuned model can be deployed with the same inference-time memory footprint as the pretrained model, as long as no few-shot exemplars are used. This highlights a key trade-off: few-shot prompting can be very effective at steering behavior but may be difficult to scale when memory is tight or when many images need to be processed in parallel.

## 4.3 Qualitative Galleries

To better understand how each configuration behaves beyond aggregate similarity metrics, I constructed a series of qualitative galleries. Each gallery focuses on a single image from the PixelProse subset and shows:

- The input image (top),
- The PixelProse reference caption, and
- The generated captions from five configurations: pretrained factual (PT + N), pretrained detailed (PT + D), pretrained detailed multimodal few-shot (PT + D\_FS), fine-tuned factual (FT + N), and fine-tuned detailed (FT + D).

These galleries make it easier to see how prompting style and LoRA fine-tuning jointly affect caption length, style, and failure modes such as hallucination and repetition.



**PixelProse (reference):**

This image displays:

A woman wearing a white and gray zip-up jacket with a hood. The jacket has two front pockets. The hood is also gray on the inside. The woman is looking to the side, and away from the camera. The woman is standing with her arms at her side. She has a neutral, perhaps slightly happy expression on her face. The background is white without any other objects or scenery. The image is a photograph of a woman posing in front of a white background, and the style is fashion-oriented.

**PT + N (pretrained, factual prompt):**

A woman wearing a grey hoodie.

**PT + D (pretrained, detailed prompt):**

A woman wearing a grey hoodie and white pants.

**PT + D\_FS (pretrained, detailed few-shot):**

This image displays a person wearing a light gray hoodie with a hood and white pants. The hoodie has a pattern of small, dark gray shapes on the front. The person is standing with their hands in their pockets, facing away from the camera. The background is plain white.

**FT + N (fine-tuned, factual prompt):**

A young woman is wearing a grey hoodie sweatshirt with a white drawstring. She is standing with her left hand in her hoodie pocket and her right hand hanging by her side. She is looking at the camera with a slight smile on her face. The background is a white wall. The image is a photograph taken from a slightly elevated angle. The woman is wearing a neutral expression on her face. The image is well lit and the colors are muted. The style of the image is casual.

**FT + D (fine-tuned, detailed prompt):**

This image displays a young woman wearing a grey hoodie sweatshirt. The sweatshirt has a hood and long sleeves. The woman is standing with her left hand in her hoodie pocket and her right hand hanging by her side. She is looking at the camera with a slight smile on her face. The background is white and there is a shadow on the left side of the woman. The image is a photograph taken from a slightly elevated angle. The woman is wearing a neutral expression on her face. The image is well lit and the colors are muted. The style of the image is casual and relaxed. *(The model begins to repeat this description later in the full output, indicating a repetition tendency after fine-tuning.)*

Figure 4: Qualitative gallery for ID 3 (woman in gray hoodie).



**PixelProse (reference):**

This image displays a stainless steel folding table. The table is made of stainless steel and has a foldable design. It is shown folded in half. The table is shown against a white background. The table is not being used, and there are no people present. The image is a photograph and is factual and terse but complete.

**PT + N (pretrained, factual prompt):**

The image shows a stainless steel industrial refrigeration unit with multiple compartments. It is designed for commercial use, likely in a food service or retail setting, to keep perishable items cold. The unit has a sturdy metal frame and a flat top surface.

**PT + D (pretrained, detailed prompt):**

The image features a stainless steel kitchen cart or cart with a drawer. The cart has a sturdy metal frame and a drawer that is open, revealing the interior of the cart. The drawer is made of stainless steel and has a handle on the side. The cart is designed for use in a kitchen or food preparation area, providing a convenient and functional storage solution.

**PT + D\_FS (pretrained, detailed few-shot):**

This image displays a stainless steel industrial freezer with multiple compartments. The freezer is designed with a sturdy metal frame and features a sliding door mechanism. The compartments are divided into sections, each with a clear, reflective surface, likely for easy visibility of the contents. The freezer appears to be set up for commercial use, possibly in a food service or retail environment. The overall design suggests functionality and durability, suitable for storing perishable goods.

**FT + N (fine-tuned, factual prompt):**

A stainless steel kitchen cart with a countertop and a drawer. The drawer is open and the countertop is empty. The cart is on a white background. The image is a photograph taken from a top-down perspective. The cart is made of metal and has a simple, utilitarian design. The countertop is made of stainless steel and has a smooth, shiny surface. The drawer is made of stainless steel and has a simple, utilitarian design. *(The caption continues by repeating this description.)*

**FT + D (fine-tuned, detailed prompt):**

This is a stainless steel kitchen cart with a countertop and a drawer. The cart is made of metal and has a silver finish. The countertop is made of stainless steel and has a smooth surface. The drawer is made of stainless steel and has a hinged door. The cart is supported by four metal legs. The cart is empty and there is nothing on the countertop. The background is white and there is no scenery. The image is a 3D render and is not a photograph. *(The model then repeats the "3D render" statement multiple times.)*

Figure 5: Qualitative gallery for ID 4 (stainless steel folding table).



**PixelProse (reference):**

This image displays: a cowboy on a horse. The cowboy has a lasso in his hand and is wearing a hat. The horse is rearing up and has one front hoof off the ground. Both the cowboy and the horse are in mid-action. There is no background scenery. The image is a line drawing in black and white.

**PT + N (pretrained, factual prompt):**

A cowboy riding a horse, depicted in a dynamic pose, with the horse rearing up on its hind legs. The cowboy is wearing a hat and traditional attire, holding the reins as the horse appears to be in motion.

**PT + D (pretrained, detailed prompt):**

The image features a black and white line drawing of a cowboy riding a horse. The cowboy is depicted in a dynamic pose, with one arm extended and the other holding the reins. The horse is in mid-gallop, with its legs lifted off the ground, suggesting movement and speed. The cowboy's attire includes a hat, boots, and a long-sleeved shirt, typical of Western wear. The horse's mane and tail are flowing, adding to the sense of motion. The background is plain, focusing attention on the cowboy and horse. There are no additional objects or details in the image, making it a simple yet evocative representation of a classic Western scene.

**PT + D\_FS (pretrained, detailed few-shot):**

The image is a black and white line drawing of a cowboy riding a horse. The cowboy is wearing a hat, a long-sleeved shirt, and jeans. He is holding the reins of the horse with one hand and has his other hand extended outward. The horse is depicted in mid-gallop, with its front legs lifted off the ground. The cowboy appears to be in motion, suggesting a sense of speed and action. The drawing is detailed, with clear lines defining the shapes and features of both the cowboy and the horse.

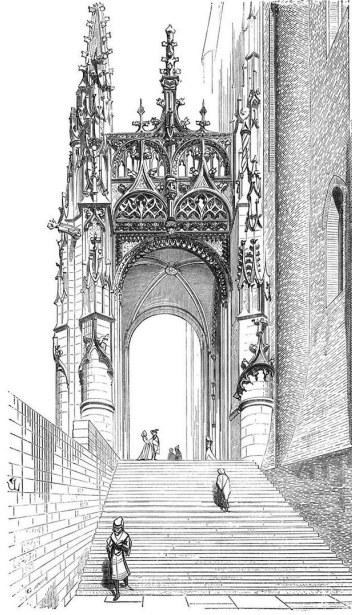
**FT + N (fine-tuned, factual prompt):**

A line drawing of a cowboy riding a horse. The cowboy is wearing a cowboy hat and a long-sleeved shirt. The horse is brown and has a saddle on its back. The cowboy is holding the reins in his right hand and is looking to the left. The horse is standing on all four legs and is in a galloping position. The background is white and there is no scenery. The drawing is in a cartoon style and is drawn in black and white. *(The caption then repeats similar sentences several times.)*

**FT + D (fine-tuned, detailed prompt):**

This is a line drawing of a cowboy riding a horse. The cowboy is wearing a cowboy hat, a long-sleeved shirt, and jeans. He is holding the reins in his right hand and is looking to the left. The horse is brown and has a saddle on its back. The horse is standing on all four legs. The background is white. The drawing is in a cartoon style. *(The remainder of the caption largely repeats this description, illustrating repetition after fine-tuning.)*

Figure 6: Qualitative gallery for ID 5 (cowboy on a horse, line drawing).



**PixelProse (reference):**

This image displays:

A black-and-white drawing of an ornate Gothic doorway. The doorway has two staircases leading up to it, and is flanked by two large buttresses. The doorway itself is topped by a gable, and is decorated with statues, pinnacles, and other intricate carvings.

**PT + N (pretrained, factual prompt):**

The image is a detailed black and white illustration of a grand Gothic-style archway, likely part of a cathedral or significant historical building. The archway features intricate carvings and ornate stonework, with pointed arches and elaborate tracery. A wide staircase leads up to the arch, with several figures depicted ascending and descending the steps. The scene captures the architectural grandeur and historical significance of the structure.

**PT + D (pretrained, detailed prompt):**

The image is a detailed black and white illustration of a grand, ornate Gothic-style archway, likely part of a cathedral or significant historical building. The archway is intricately decorated with pointed arches, pinnacles, and flying buttresses, characteristic of Gothic architecture. The scene is set on a wide staircase leading up to the archway, with several figures depicted in various poses, suggesting movement and activity. The figures appear to be dressed in period attire, indicating a historical setting. The overall composition conveys a sense of grandeur and historical significance, with attention to architectural detail and the depiction of human figures in a dynamic scene.

**PT + D-<sub>FS</sub> (pretrained, detailed few-shot):**

The image depicts a detailed black and white illustration of a grand, ornate Gothic-style archway. The arch is richly decorated with intricate carvings and pointed spires, characteristic of Gothic architecture. The structure features a large central arch flanked by two smaller arches, all adorned with elaborate tracery and statues. The steps leading up to the arch are wide and steep, with a few figures ascending and descending them. The figures are dressed in period clothing, suggesting a historical or fantasy setting. The overall composition conveys a sense of grandeur and historical significance.

**FT + N (fine-tuned, factual prompt):**

A black and white line drawing of a Gothic cathedral with a large entrance and a set of stairs leading up to it. There are three people on the stairs, two of them are standing and one is sitting on the steps. The cathedral is made of stone and has a large tower on the left side. The entrance to the cathedral is made of stone and has a large archway decorated with intricate carvings. *(Later in the full output, the model emphasizes repeatedly that the image is a drawing rather than a photograph.)*

**FT + D (fine-tuned, detailed prompt):**

This is a black and white line drawing of a Gothic cathedral. The cathedral has a large central nave with a pointed arch and a large rose window at the top. There are two towers on either side of the nave, each with a pointed roof and a large spire on top. The walls of the cathedral are made of stone and there are many small windows along the walls. There is a large staircase leading up to the cathedral from the bottom of the image, with three people on the stairs. *(The caption later repeats that the image is a drawing and not a photograph.)*

Figure 7: Qualitative gallery for ID 6 (ornate Gothic doorway drawing).



**PixelProse (reference):**

This image displays: a field of tall, yellow grass in the foreground. There is a large hill covered in green trees in the background. A rainbow can be seen in the sky above the hill. The rainbow is a full rainbow, with the red end on the left and the purple end on the right. There are no people or other objects visible in the image. The image is bright and sunny, and it appears to have been taken from a high angle. The answer to the question in the original caption is: yes, your sister did find the pot of gold, or at least the end of the rainbow.

**PT + N (pretrained, factual prompt):**

A rainbow is visible in the sky above a grassy field with trees in the background.

**PT + D (pretrained, detailed prompt):**

The image captures a serene rural landscape with a vibrant rainbow arching across the sky. The rainbow, with its distinct colors, adds a touch of magic to the scene. The foreground features a field, possibly a farm or a pasture, with a fence running along its edge. Beyond the field, there's a line of trees and bushes, providing a natural boundary to the area. The background is dominated by a hill or mountain, covered in lush greenery, suggesting a location that is rich in natural beauty. The overall atmosphere of the image is peaceful and idyllic, evoking a sense of tranquility and connection with nature.

**PT + D.FS (pretrained, detailed few-shot):**

The image depicts a serene landscape featuring a vibrant rainbow arching across the sky. The rainbow displays a spectrum of colors, including red, orange, yellow, green, blue, indigo, and violet, creating a striking contrast against the clear sky. Below the rainbow, there is a lush, green field with tall grasses and a variety of trees and bushes. The trees are dense and appear to be in full foliage, suggesting a healthy, thriving environment. The background consists of a hill or mountain covered with dense forest, adding depth to the scene. The overall atmosphere of the image is peaceful and picturesque, evoking a sense of tranquility and natural beauty.

**FT + N (fine-tuned, factual prompt):**

A rainbow is seen in the sky above a green forest. The rainbow is a spectrum of colors, with red on the left and violet on the right. The colors are orange, yellow, green, blue, indigo, and violet. The rainbow is a natural phenomenon that occurs when sunlight is refracted by water droplets in the air. The rainbow is a beautiful and awe-inspiring sight that can be seen by anyone who is willing to look. It is a reminder of the beauty and wonder of the natural world. *(The caption continues by repeating a generic explanation of rainbows.)*

**FT + D (fine-tuned, detailed prompt):**

This image shows a rainbow in the sky. The rainbow is a spectrum of colors, with red on the left and violet on the right. The colors are orange, yellow, green, blue, indigo, and violet. The rainbow is in the sky, and there is a hill in the background. The hill is covered in trees. The sky is blue. The image is a photograph. *(Later in the caption, the model repeats the explanation about the rainbow and viewpoint several times.)*

Figure 8: Qualitative gallery for ID 7 (rainbow over hill).



**PixelProse (reference):**

This image displays:

A woman wearing a white and gray zip-up jacket with a hood. The jacket has two front pockets. The hood is also gray on the inside. The woman is looking to the side, and away from the camera. The woman is standing with her arms at her side. She has a neutral, perhaps slightly happy expression on her face. The background is white without any other objects or scenery. The image is a photograph of a woman posing in front of a white background, and the style is fashion-oriented.

**FT + D\_FS (fine-tuned, detailed few-shot, failure case):**

*This image displays a pink background with black text. The text reads: "I just want to lose weight while staying in bed, watching TV, and eating Girl Scout cookies. Is that really too much to ask?" The text is in a white speech bubble on the left side of the image. On the right side of the image is a cartoon drawing of a person sitting on a couch with polka dots. The person is wearing pajamas and has their head in their hands. There is a box of Girl Scout cookies on the floor in front of the couch. The image is drawn in a simple cartoon style.*

**Explanation of failure:** Instead of captioning the query image (a woman wearing a gray jacket), the FT + D\_FS model reproduced the caption of *Exemplar 3* almost verbatim. This demonstrates the primary pathology observed when combining LoRA fine-tuning with multimodal few-shot prompting:

- The model overfits to the fine-tuned single-image prompt format, and
- When multiple exemplar images are introduced, it sometimes "locks on" to an exemplar and ignores the query image entirely.

This is one of the clearest examples of cross-image confusion caused by mixing LoRA-adapted layers with a multi-image context structure not seen during training.

Figure 9: Failure case for ID 3: the FT + D\_FS model confuses an exemplar with the query image and outputs the exemplar’s caption verbatim.

## 5 Discussion

The experiments suggest several key observations about the relative strengths and weaknesses of LoRA fine-tuning and multimodal few-shot prompting for dense captioning.

**Effectiveness of Few-Shot Prompting.** For the pretrained model, multimodal few-shot prompting with three exemplars consistently improves automatic metrics and produces captions that are qualitatively closer to PixelProse outputs. The exemplars seem to serve two roles: they demonstrate the desired level of detail, and they implicitly bias the model toward the PixelProse style (including sentence structure and lexical choices). This indicates that few-shot prompting is a viable way to adapt general-purpose VLMs to dense captioning without any weight updates.

**Limited Gains from Small-Scale LoRA Fine-Tuning.** Due to resource constraints and the inability to run stable multi-GPU training, the LoRA fine-tuning in this project was limited to a relatively small subset of PixelProse and a modest number of optimization steps. Under these conditions, the fine-tuned model shows only small improvements over the pretrained model when using the same prompts. In particular, FT + D is competitive with PT + D but does not dramatically outperform it. This suggests that when training data and compute are limited, prompt engineering and few-shot prompting may yield comparable or larger gains than a small LoRA run.

**Failure of the Combined FT + Few-Shot Setting.** The most surprising result is that combining LoRA fine-tuning with detailed few-shot prompting (FT + D.FS) significantly *degrades* performance. The generated captions are often extremely long, repetitive, or misaligned with the query image. In several examples, the model appears to describe one of the exemplar images instead of the query, or it loops over phrases already present in the exemplars.

A plausible explanation is that the LoRA adapter was trained solely on single-image inputs paired with the detailed prompt, without any few-shot structure. During training, the model never sees multiple images in the same sequence or captions that refer to “exemplar” and “query” images. At evaluation time, the FT + D.FS configuration suddenly introduces multiple images and a more complex prompt. The adapted language model may therefore overfit to patterns in the exemplar captions and struggle to disentangle which visual features correspond to the current generation target.

This effect underscores an important lesson: fine-tuning and prompting must be aligned. If a model is fine-tuned under one prompt distribution and deployed under a different prompt distribution (for example, adding few-shot examples or changing instruction phrasing), can degrade performance on the dense captioning task.

**Memory Cost of Multimodal Few-Shot Prompting.** Another important observation is the memory footprint of multimodal few-shot prompting. When using three exemplar images plus the query, the system consumed roughly 50 GB of VRAM across my GPUs. In settings where GPU memory is limited or where batched inference over many images is required, this overhead may be prohibitive. In contrast, once a LoRA adapter is trained, it has essentially no additional inference-time memory cost beyond the base model.

This trade-off suggests that few-shot prompting is especially attractive in scenarios with ample memory and relatively low throughput requirements (for example, interactive captioning tools), while fine-tuning may be more appropriate when serving large numbers of requests or running on edge devices.

**Implications for Future Systems.** Taken together, these results argue for more holistic design of captioning systems. Few-shot prompting is a powerful tool, but aligning the training and inference prompts is crucial. When fine-tuning is used, it should ideally be performed with the same prompt structure (including few-shot exemplars) that will be used at deployment. Conversely, when only prompting is feasible, careful selection of exemplars and anti-repetition decoding strategies can yield strong gains without any training.

## 6 Conclusion

This project set out to compare LoRA fine-tuning and multimodal few-shot prompting as strategies for improving detailed image captioning with InternVL 3.5 2B on the PixelProse dataset. Within the constraints of a single-semester project and single-GPU training, I found that:

- Multimodal few-shot prompting with three exemplar image-caption pairs substantially improves similarity to PixelProse captions for the pretrained model, both quantitatively and qualitatively.
- The small-scale LoRA fine-tuning run used here provides only modest additional gains when evaluated under the same prompts, likely due to limited training data and compute.
- Combining LoRA fine-tuning with a prompt structure that is very different from the training prompt (detailed few-shot) can lead to worse performance and pathological behavior such as repetition and off-target descriptions.
- Few-shot prompting carries a significant VRAM cost when multiple images are processed in a single context, which may limit its applicability in memory-constrained settings.

From a learning perspective, I gained practical experience setting up a modern multimodal pipeline, implementing LoRA fine-tuning for a large vision-language model, and designing careful evaluations that compare prompting and training-based approaches. The project also reinforced the importance of aligning training and inference conditions and of combining automatic metrics with qualitative analysis.

Future work could pursue several directions: scaling LoRA training to multi-GPU settings and larger subsets of PixelProse, explicitly training the model on few-shot prompts that mirror the evaluation setup, experimenting with anti-repetition loss functions or decoding strategies, and exploring more diverse exemplar selection strategies to reduce overfitting to a small set of examples.

## My Contribution

I worked alone on this project, so all aspects of the work were my responsibility. Specifically, I:

- Implemented the PixelProse subset construction script, downloaded and filtered the data, and organized the dataset into train, validation, and evaluation splits.
- Set up the InternVL 3.5 2B environment, including image tiling, preprocessing, and the captioning interface.
- Designed and implemented the LoRA fine-tuning pipeline on the language model component of InternVL, including dataset loading, batching, multimodal input construction, and checkpoint saving.
- Implemented the evaluation script for running multiple prompting configurations (neutral, detailed, and detailed few-shot) on both pretrained and fine-tuned models, and ensured exemplar IDs were excluded from the evaluation pool.
- Computed and analyzed automatic metrics (BLEU, METEOR, CIDEr, BERTScore), created summary CSV files, and constructed qualitative galleries of PixelProse and generated captions.
- Performed the analysis and discussion of results, including the identification of failure modes when combining LoRA fine-tuning with few-shot prompting and the assessment of VRAM usage during multimodal prompting.
- Wrote this report and organized the figures, tables, and references.

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022.
- [2] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
- [3] Vasu Singla, Kaiyu Yue, Sukriti Paul, Reza Shirkavand, Mayuka Jayawardhana, Alireza Ganjdanesh, Heng Huang, Abhinav Bhatele, Gowthami Somepalli, and Tom Goldstein. From pixels to prose: A large dataset of dense image captions, 2024.
- [4] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, Zhaokai Wang, Zhe Chen, Hongjie Zhang, Ganlin Yang, Haomin Wang, Qi Wei, Jinhui Yin, Wenhao Li, Erfei Cui, Guanzhou Chen, Zichen Ding, Changyao Tian, Zhenyu Wu, Jingjing Xie, Zehao Li, Bowen Yang, Yuchen Duan, Xuehui Wang, Zhi Hou, Haoran Hao, Tianyi Zhang, Songze Li, Xiangyu Zhao, Haodong Duan, Nianchen Deng, Bin Fu, Yinan He, Yi Wang, Conghui He, Botian Shi, Junjun He, Yingdong Xiong, Han Lv, Lijun Wu, Wenqi Shao, Kaipeng Zhang, Huipeng Deng, Biqing Qi, Jiaye Ge, Qipeng Guo, Wenwei Zhang, Songyang Zhang, Maosong Cao, Junyao Lin, Kexian Tang, Jianfei Gao, Haiyan Huang, Yuzhe Gu, Chengqi Lyu, Huanze Tang, Rui Wang, Haijun Lv, Wanli Ouyang, Limin Wang, Min Dou, Xizhou Zhu, Tong Lu, Dahua Lin, Jifeng Dai, Weijie Su, Bowen Zhou, Kai Chen, Yu Qiao, Wenhao Wang, and Gen Luo. Internvl3.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency, 2025.