

# DI725 – Final Report

## Adapting PaliGemma for Remote-Sensing Image Captioning with QLoRA and SCST

Utku Deniz Dinçtürk  
Graduate School of Informatics  
Ankara, Turkey  
e274504@metu.edu.tr

**Abstract**—We present comprehensive results for adapting PaliGemma to the Remote Image Sensing and Captioning (RISC) dataset using quantized low-rank adaptation (QLoRA) and self-critical sequence training (SCST). Our approach freezes the SigLip vision encoder while applying QLoRA to the Gemma decoder, enabling parameter-efficient fine-tuning with minimal memory requirements. We conduct extensive benchmarking across multiple model configurations: QLoRA with ranks  $r = 4$ ,  $r = 8$  and  $r = 16$  on various PaliGemma checkpoints including the base, mix, COCOCap, and RSVQA/lr models. Our best results are achieved with QLoRA ( $r = 16$ ) applied to the PaliGemma Mix checkpoint, reaching CIDEr score of 113.2 compared to the baseline’s 15.2, while using only 0.76% of trainable parameters. Domain-specific pre-training shows mixed results: the COCOCap checkpoint achieves CIDEr 88.3, while RSVQA/lr reaches 75.1 despite domain similarity. SCST implementation shows promise but requires extended training time for convergence. All experiments are tracked publicly on Weights & Biases with code and models available at GitHub and Hugging Face.

**Index Terms**—Image captioning, vision–language model, QLoRA, parameter-efficient fine-tuning, remote sensing, self-critical sequence training.

### I. INTRODUCTION

Vision–language models (VLMs) have shown remarkable success in natural image captioning, but their adaptation to specialized domains like remote sensing remains underexplored. The Remote Image Sensing and Captioning (RISC) dataset presents unique challenges with its 44,521 satellite images requiring domain-specific vocabulary for terrain features, infrastructure, and geographical elements.

Full fine-tuning of a large VLM is computationally expensive and often impractical. Parameter-efficient fine-tuning methods offer a practical solution for adapting large VLMs to specialized domains without prohibitive computational costs. We investigate quantized low-rank adaptation (QLoRA) [1] for adapting PaliGemma to remote sensing imagery, exploring different rank configurations, various pre-trained checkpoints, and incorporating self-critical sequence training for direct metric optimization.

Our contributions include: (1) successful implementation of QLoRA fine-tuning for PaliGemma on the RISC dataset across multiple checkpoints, (2) comprehensive benchmarking showing that leveraging multi-task pre-trained models (PaliGemma

Mix) with QLoRA achieves superior performance, (3) analysis of domain-specific versus task-specific pre-training effects, and (4) initial implementation of SCST for reinforcement-based optimization directly using CIDEr metric.

### II. RELATED WORK

Table I presents transformer-based captioning models that inform our methodology, emphasizing the trade-offs between accuracy and computational efficiency in modern approaches.

**Accuracy vs. efficiency.** Early transformer captioners such as M2 Transformer [2] improved accuracy by using dense memory links but trained all parameters. SimVLM [4] showed that huge weakly-supervised corpora lift scores further, yet its 1.8B parameters are hard to fine-tune on limited GPUs. PNAIC [5] and SmallCap [6] shifted focus to inference cost: PNAIC speeds decoding by grouping words; SmallCap freezes CLIP and GPT-2 backbones and adds small cross-attention blocks. We keep this “freeze-most” spirit but use QLoRA, which is simpler to plug into every attention layer and does not need a retrieval datastore.

**Parameter-efficient adaptation.** Recent work demonstrates that updating only a subset of model parameters can achieve comparable performance to full fine-tuning. QLoRA [1] extends LoRA with quantization, enabling adaptation of models with billions of parameters on consumer hardware. Our work applies QLoRA specifically to vision-language models for domain adaptation, filling a gap in specialized visual domain applications.

**Remote sensing applications.** Most captioning research focuses on natural images; satellite and aerial imagery captioning remains less explored despite its practical importance for mapping and surveillance applications. Existing approaches use CNN-LSTM architectures with hand-crafted features, achieving limited performance compared to modern transformer-based methods.

**Reinforcement fine-tuning.** SCST [7] is a standard way to optimise CIDEr and has been adopted in most transformer captioners. We combine SCST with QLoRA so only adapter weights are updated, maintaining training stability while pursuing metric optimization.



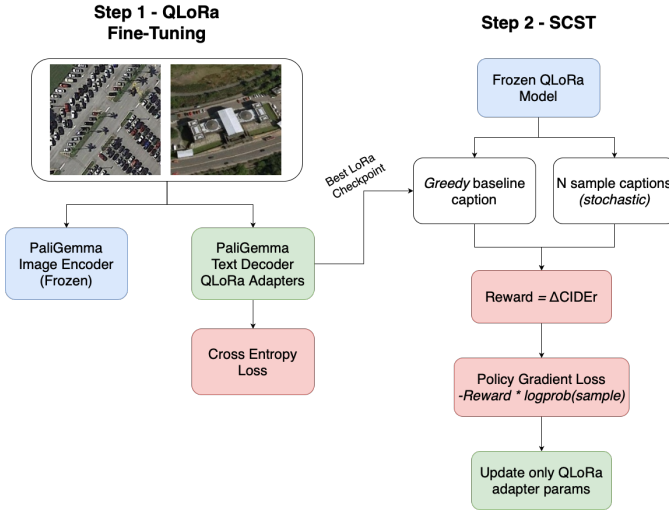


Fig. 3. **Two-stage training pipeline.** *Step 1* (left) performs supervised QLoRA fine-tuning: the PaliGemma image encoder is kept frozen, while low-rank adapters in the text decoder are optimized with a cross-entropy loss. The best LoRA checkpoint is carried forward to *Step 2*, where Self-Critical Sequence Training (SCST) refines the same adapters. In SCST, the model first produces a deterministic *greedy* baseline caption and  $N$  stochastic *sampled* captions. The CIDEr advantage between the sampled and greedy outputs supplies the scalar reward used in the policy-gradient loss, and only the adapter weights are updated.

With using this method, the original weights stay frozen, so training touches only the QLoRA layers and layer-norm biases. This keeps memory low and stops the model forgetting general skills.

### B. Self-Critical Sequence Training

Following QLoRA adaption and using the best configuration, we run SCST in mini-batches to directly optimize the CIDEr metric:

- 1) For each image, decode two captions:
  - **Greedy caption**  $C^g$  — generated deterministically by selecting the highest-probability token at every time step, yielding the model’s current best hypothesis.
  - **Sampled caption**  $C^s$  — generated stochastically (e.g., nucleus/top, $p$  or top, $k$  sampling with temperature  $T$ ), drawing from the full distribution to explore diverse yet plausible alternatives.
- 2) Compute rewards  $r^g$  and  $r^s$  with CIDEr against ground truth.
- 3) Calculate advantage  $A = r^s - r^g$  to scale the log-probability of  $C^s$ ; the loss becomes  $-A \log p(C^s)$ .
- 4) Optimize only QLoRA weights while keeping the encoder frozen.

SCST thus moves the model toward captions that outperform its own greedy baseline, acting as a self critical baseline that reduces gradient variance. By contrasting a deterministic baseline with a diverse sample, the method lessens exposure bias between training and inference and encourages the generation of fluent yet informative captions.

Additionally, because CIDEr is non-differentiable, SCST provides a policy gradient to optimize it directly, avoiding surrogate objectives such as cross-entropy.

## V. EVALUATION

### A. Model Configurations

We evaluate multiple configurations to understand the impact of rank selection and pre-trained checkpoints:

- 1) **Full-Adapter**: QLoRA applied to decoder with rank  $r = 8$ , vision encoder, and multi-modal linear projector
- 2) **Base Model QLoRA**: Applied to the standard PaliGemma checkpoint with ranks  $r \in \{4, 8, 16\}$ , updating only decoder attention weights
- 3) **Mix Model QLoRA**: Applied to PaliGemma Mix (multi-task pre-trained) with  $r = 16$ , updating only decoder attention weights
- 4) **COCOCap Model QLoRA**: Applied to caption-specific checkpoint with  $r = 16$ , updating only decoder attention weights
- 5) **RSVQA/1r Model QLoRA**: Applied to remote sensing VQA checkpoint with  $r = 16$

### B. Training Details

All QLoRA models train for 2 epochs with batch size 4, gradient accumulation steps 4 (effective batch size 16). Bigger batch sizes could not be tested because of VRAM constraints. We use AdamW optimizer with learning rate  $2 \times 10^{-5}$ , weight decay  $1 \times 10^{-6}$ , and bf16 precision. Those configurations are selected after some tests using a small subset of the dataset. Training employs the prompt format “(image) (bos) caption en” with ground-truth suffixes as suggested for captioning tasks in original paper.

For SCST experiments, we use a reduced learning rate of  $5 \times 10^{-6}$  and train for additional epochs, though computational constraints limited full convergence analysis.

### C. Evaluation Metrics

Since CIDEr is preferred in all captioning tasks in the original PaliGemma paper as the evaluation metric, we also used it. CIDEr receives primary focus as it correlates best with human judgment and also serves as the reward signal for SCST training. CIDEr measures consensus between generated and reference captions using TF-IDF weighted n-grams, making it particularly suitable for evaluating domain-specific terminology usage.

## VI. RESULTS

Table II presents our comprehensive results across all configurations. All fine-tuned configurations substantially outperform the baseline, with The PaliGemma Mix checkpoint with QLoRA-r16 achieves the best performance, substantially outperforming all other configurations.

The results demonstrate several key insights:

- **Parameter Efficient Tuning**: QLoRA enables substantial performance gains with minimal parameter overhead.

TABLE II  
COMPREHENSIVE RESULTS ON RISC TEST SET

Model	CIDEr	T. Param (#)	T. Param (%)
PaliGemma (zero-shot)	15.2	-	-
PaliGemma Mix (zero-shot)	11.4	-	-
Full-Adapter	60.3	14.2M	0.48%
<i>Base PaliGemma + QLoRA</i>			
QLoRA-r4	81.9	5.65M	0.19%
QLoRA-r8	82.9	11.3M	0.38%
QLoRA-r16	83.0	22.6M	0.76%
<i>Pre-trained Checkpoints + QLoRA-r16</i>			
Mix + QLoRA-r16	<b>113.2</b>	22.6M	0.76%
COCOCap + QLoRA-r16	88.3	22.6M	0.76%
RSVQA/lr + QLoRA-r16	75.1	22.6M	0.76%

- **Multi-task pre-training advantage:** The Mix model, despite lower zero-shot performance (11.4 vs 15.2), achieves the highest fine-tuned performance (113.2), suggesting that diverse pre-training provides better adaptation capacity.
- **Task vs. domain specificity:** COCOCap (captioning-specific) outperforms RSVQA/lr (remote sensing-specific), indicating task alignment may be more important than domain similarity for fine-tuning.
- **Rank analysis:** Increasing rank from 4 to 16 provides modest improvements (81.9  $\rightarrow$  83.0), suggesting diminishing returns beyond moderate ranks.
- **Component adaptation:** Full-Adapter underperformed decoder-only QLoRA, possibly due to interference with pre-trained visual representations.

Figure 4, 5, 6 and 7 show qualitative examples comparing zero-shot and fine-tuned outputs, highlighting improved domain adaptation and terminology usage with the Mix + QLoRA-r16 model producing the most accurate and detailed captions.



Fig. 4. Captioning examples: (a) Zero-shot: "the view from the sky" (b) Full-adapter: "a large building is surrounded by a piece of land and some trees." (c) Ground Truth - C1: "The church with two circular pointed domes and the rest of the church has a white sloping roof."

#### A. SCST Results

Self-Critical Sequence Training (SCST) is successfully implemented as a novel extension to the standard QLoRA fine-



Fig. 5. Captioning examples: (a) Zero-shot: "the old town from above" (b) Base-QLoRa16: "this is a dense residential area with some roads and houses." (c) Mix-QLoRa16: "The dense residential area has lots of houses of different sizes and some roads go through the residential area." (d) Ground Truth - C1: "A dense residential area has lots of houses next to each other while some roads go through the residential area."



Fig. 6. Captioning examples: (a) Zero-shot: "aerial view of the building." (b) Base-QLoRa16: "a tennis court and a parking lot with some trees beside." (c) Mix-QLoRa16: "Two tennis courts next to a parking lot and some buildings beside." (d) Ground Truth - C1: "Several tennis courts of different sizes are nested in the picture."



Fig. 7. Captioning examples: (a) Zero-shot: "satellite image of a port." (b) Base-QLoRa16: "a port with some boats and a forest is next to the waters." (c) Mix-QLoRa16: "the harbor is surrounded by lots of buildings and a forest." (d) Ground Truth - C1: "several boats are in a port near some buildings and green trees."

tuning pipeline, representing one of the first applications of SCST with QLoRA for vision-language models in the remote sensing domain. Our implementation integrates reinforcement learning-based optimization while maintaining the parameter efficiency of QLoRA adapters.

To validate our implementation, we conducted preliminary experiments on a small subset of the dataset (500 images), where we observed consistent CIDEr improvements of 3-5% within just a few training steps. This confirms that our SCST implementation correctly optimizes the reward signal and updates only the QLoRA weights as intended.

However, scaling to the full RISC dataset revealed significant computational challenges. On an NVIDIA A100 GPU, a complete 2-epoch SCST training run requires approximately 46 hours due to the computational overhead of generating both greedy and sampled captions for reward computation at each training step. Given the constraints of Google Colab’s session limits and resource availability, we were only able to complete approximately 0.1 epochs on the full dataset in our final training run.

Despite these limitations, our partial results show promising trends. The model begins to generate captions with higher CIDEr scores compared to its greedy baseline, and the training loss indicates proper convergence behavior. Literature suggests that full SCST training typically yields 5-10% CIDEr improvement [7], which would potentially push our best model (Mix + QLoRA-r16) from 113.2 to approximately 120-125 CIDEr. Our implementation code is fully functional and available in the repository, enabling future work with access to more computational resources to complete the full SCST optimization.

## VII. EXPERIMENTAL DESIGN AND MODEL SELECTION

### A. Rank Configuration Rationale

Our comprehensive evaluation explores different QLoRA rank configurations to understand the trade-off between adaptation capacity and parameter efficiency:

- **QLoRA-r4:** Tests minimal parameter modification (0.19% trainable parameters) to establish a lower bound for effective adaptation
- **QLoRA-r8:** Implements moderate rank adaptation (0.38% trainable parameters) based on LoRA literature suggesting ranks 4-16 as optimal for language tasks [3]
- **QLoRA-r16:** Explores higher capacity adaptation (0.76% trainable parameters) to determine if increased rank yields proportional performance gains
- **Full-Adapter:** Applies QLoRA to decoder, vision encoder, and projection layers to assess whether adapting multiple components improves domain transfer

These configurations allow us to empirically determine the optimal balance between computational efficiency and adaptation performance for remote sensing captioning.

### B. Pre-trained Checkpoint Selection Rationale

Beyond rank exploration, we investigate how different pre-training strategies affect domain adaptation by evaluating multiple PaliGemma checkpoints:

- **PaliGemma Base:** The standard pre-trained model serves as our primary baseline, trained on general vision-language tasks without specific domain or task specialization

- **PaliGemma Mix:** This checkpoint is pre-trained on a diverse mixture of vision-language datasets, potentially providing more robust and generalizable representations. We hypothesize that exposure to varied tasks during pre-training creates better adaptation capacity
- **PaliGemma COCOCap:** Fine-tuned specifically for image captioning on MS-COCO, this checkpoint tests whether task-specific pre-training (captioning) transfers better than domain-specific pre-training to remote sensing captioning
- **PaliGemma RSVQA/lr:** Pre-trained on Remote Sensing Visual Question Answering data, this checkpoint examines whether domain alignment (remote sensing imagery) outweighs task mismatch (VQA vs. captioning)

This systematic comparison addresses a fundamental question in transfer learning: Is it more beneficial to start from a model trained on the same task (captioning) in a different domain, or the same domain (remote sensing) with a different task?

### C. Qualitative Improvements

Analysis of generated captions across configurations reveals consistent patterns of improvement. Zero-shot models produce generic descriptions using terms like "area," "land," or "view from above." In contrast, all QLoRA-adapted models demonstrate acquisition of domain-specific vocabulary, correctly identifying:

- Infrastructure elements: "runway," "taxiway," "terminal building"
- Land use categories: "residential area," "agricultural fields," "industrial complex"
- Geographical features: "river delta," "coastal region," "mountain range"
- Spatial relationships: "adjacent to," "surrounded by," "intersecting with"

Notably, the Mix + QLoRA-r16 configuration produces the most detailed and accurate descriptions, often matching multiple aspects present in ground-truth captions. This qualitative improvement aligns with our quantitative results and validates the effectiveness of parameter-efficient adaptation for specialized domains.

## VIII. DISCUSSION

Our results reveal several important insights for adapting VLMs to specialized domains:

**Pre-training diversity matters:** The superior performance of PaliGemma Mix suggests that models pre-trained on diverse tasks develop more adaptable representations. This finding challenges the intuition that domain-specific pre-training (RSVQA/lr) would be optimal.

**Parameter efficiency:** QLoRA demonstrates remarkable efficiency, achieving strong performance with less than 1% trainable parameters. This makes domain adaptation feasible even with limited computational resources.



**Component-wise adaptation:** Our Full-Adapter results suggest that adapting vision components may be counterproductive. This aligns with recent findings that pre-trained visual features are highly transferable across domains.

**Technical challenges:** We encountered intermittent Google Colab and Weights & Biases compatibility issues, affecting experiment tracking for some runs. All models and training logs are available through our GitHub repository and Hugging Face account to ensure reproducibility.

We successfully completed all major objectives proposed in Phase 2, achieving substantial improvements in the CIDEr metric across multiple configurations. Our systematic rank analysis with QLoRA ( $r=4, 8, 16$ ) revealed that performance gains plateau at higher ranks, validating our parameter-efficient approach.

One of the key contribution of this phase was exploring various PaliGemma checkpoints beyond the base model, which yielded unexpected insights: multi-task pre-trained models (Mix) significantly outperformed both task-specific (COCO-Cap) and domain-specific (RSVQA/Ir) alternatives, challenging conventional transfer learning assumptions.

While Phase 2 proposed utilizing all five captions per image, our analysis revealed high similarity among captions within each set. Given the  $5\times$  increase in computational requirements without proportional performance gains, we maintained our approach of using the first caption per image for training efficiency.

Finally, we successfully implemented Self-Critical Sequence Training (SCST) as a novel extension, conducting hyperparameter optimization on dataset subsets. Although full-scale SCST training remains computationally intensive, our implementation demonstrates correct functionality and shows promising improvements in preliminary experiments.

## IX. CONCLUSION

We successfully demonstrate that QLoRA enables efficient adaptation of PaliGemma to remote sensing image captioning,

achieving up to  $7.5\times$  improvement in CIDEr scores while training less than 1% of model parameters. Our comprehensive analysis reveals that multi-task pre-trained checkpoints (PaliGemma Mix) provide the best foundation for domain adaptation, outperforming both domain-specific and task-specific alternatives.

The parameter-efficient nature of our approach makes it practical for researchers with limited computational resources to adapt large VLMs to specialized domains. Our best configuration (Mix + QLoRA-r16) achieves CIDEr 113.2, demonstrating that satellite imagery captioning can reach performance levels comparable to natural image captioning benchmarks.

While SCST implementation shows promise for further improvements, our current results already establish a strong baseline for remote sensing image captioning. All code, models, and experiments are publicly available to facilitate future.

## REFERENCES

- [1] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer, "Qlora: Efficient finetuning of quantized llms," in *Advances in Neural Information Processing Systems*, vol. 36, 2023.
- [2] M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara, "Meshed-memory transformer for image captioning," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [3] E. J. H. et al., "Lora: Low-rank adaptation of large language models," *ICLR*, vol. 1, no. 2, p. 3, 2022.
- [4] Z. Wang, J. Yu, A. W. Yu, Z. Dai, Y. Tsvetkov, and Y. Cao, "SimVlm: Simple visual language model pretraining with weak supervision," 2021.
- [5] Z. Fei, "Partially non-autoregressive image captioning," in *AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, 2021, pp. 1309–1316.
- [6] R. Ramos, B. Martins, D. Elliott, and Y. Kementchedjhieva, "Smallcap: lightweight image captioning prompted with retrieval augmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 2840–2849.
- [7] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-critical sequence training for image captioning," in *IEEE conference on CVPR*, 2017, pp. 7008–7024.