



École Supérieure de Sidi Bel Abbès (ESI-SBA)

AISD

Few-Shot Remote Sensing Image Scene Classification

Using Vision-Language Models: A Comparative Study

Groupe :

Baidar Samir
Ali Abbou Oussama
Djeziri Oussama
Senhadji M. Said

Encadré par :

Chaib Souleyman

June 10, 2025

TABLE DES MATIÈRES
CONTENTS

Table des Matières	1
I Introduction - Overview of RSI Scene Classification	2
II Motivation	2
III Problematic	3
IV Related Work	3
IV-A Classical Few-Shot Learning Approaches	3
IV-B Vision-Language Model Approaches	3
References	3
V Proposed Approach	3
V-A Pipeline Outline	3
V-B Model Zoo	4
V-B1 CNN (AlexNet)	4
V-B2 CLIP ViT-B/16	4
V-B3 CLIP ViT-L/14@336px	4
V-B4 BLIP-2 Image Encoder	4
V-B5 Remote-CLIP	5
V-C Prototype Classifier	5
VI Datasets	5
VI-A NWPU-RESISC45	5
VI-B UC-Merced (UCM)	6
VI-C AID-30	6
VI-D MLRS-Net	6
VII Evaluation Protocol and Metrics	6
VII-A Episodic Few-Shot Protocol	6
VII-B Primary Metric	7
VII-C Complementary Metrics	7
VII-D Statistical Significance Tests	7
VIII Comparative Table	7
IX Conclusion and Future Work	11
References	11

Few-Shot Remote Sensing Image Scene Classification Using Vision-Language Models: A Comparative Study

Abstract—Remote Sensing Image Scene Classification (RSI-SC) plays a vital role in geospatial analysis, environmental monitoring, and land-use planning. With the increasing availability of high-resolution satellite imagery, automated scene classification has become essential. While Convolutional Neural Networks (CNNs) have shown promising results in this domain, their reliance on large-scale labeled data limits their generalizability to new scenes or domains with few examples. In this final year project, we aim to investigate the effectiveness of recent Vision-Language Models (VLMs) and other deep learning approaches under few-shot learning settings for RSI-SC.

We conduct a comparative study using four diverse remote sensing datasets—NWPU-RESISC45, AID, UC Merced, and MLRSNet—by implementing five different model architectures inspired by research literature, including CNN, Vision Transformer (ViT), CLIP, BLIP-2, and a hybrid transformer-based method. Our evaluation follows the standard N-way K-shot protocol with 1-shot and 5-shot settings to reflect real-world label-scarce environments. This work not only provides insights into the strengths and weaknesses of each approach but also serves as a hands-on experience in reading, implementing, and critically analyzing state-of-the-art research in the field of remote sensing and machine learning.

Keywords: few-shot learning, remote sensing, scene classification, deep learning, vision-language models, CLIP, BLIP-2, CNN, prototype learning, ViT, VLMs

I. INTRODUCTION - OVERVIEW OF RSI SCENE CLASSIFICATION

Remote sensing image scene classification (RSI-SC) is a fundamental task in remote sensing applications, involving the categorization of aerial and satellite images into predefined land-use and land-cover classes. Traditional approaches have achieved excellent performance when abundant labeled training data is available, with state-of-the-art methods reaching 95-99% top-1 accuracy on standard benchmarks.

The task involves several unique challenges compared to natural image classification. Remote sensing images exhibit high intra-class variance due to different imaging conditions, seasonal variations, and geographic diversity. Additionally, inter-class similarity can be high, particularly between urban categories or natural landscapes. The overhead perspective introduces spatial relationships and geometric patterns that are distinct from ground-level photography.

Recent advances have been driven by deep learning architectures, particularly convolutional neural networks (CNNs) and more recently, vision transformers (ViTs). These models have demonstrated superior feature extraction capabilities for capturing both local texture details and global spatial

arrangements in remote sensing imagery. However, the success of these approaches heavily depends on the availability of large-scale annotated datasets, which are expensive and time-consuming to create for remote sensing applications.

II. MOTIVATION

Remote Sensing Image Scene Classification (RSI-SC) is a critical task in various real-world applications such as urban planning, environmental monitoring, disaster response, agricultural assessment, and military surveillance. The ability to automatically and accurately classify scenes from satellite or aerial imagery enables timely decision-making and efficient resource allocation across these domains. As satellite data becomes increasingly available and high in resolution, the demand for intelligent and scalable scene classification solutions continues to grow.

However, traditional deep learning approaches especially Convolutional Neural Networks (CNNs) rely heavily on large-scale annotated datasets to achieve high performance. In the context of remote sensing, labeling such datasets is not only time-consuming and costly, but also requires domain-specific expertise, making large annotated datasets difficult to obtain. This bottleneck significantly limits the generalizability of traditional models to new or rare geographic regions, novel land-use patterns, and different environmental conditions.

To address these limitations, our project focuses on **Few-Shot Learning (FSL)** approaches, which aim to train models capable of generalizing from a limited number of labeled examples per class. Few-shot learning is especially important in real-world remote sensing scenarios, where rapid adaptation to new scenes or datasets with minimal supervision is crucial.

In recent years, **Vision-Language Models (VLMs)** such as CLIP and BLIP have demonstrated remarkable few-shot and even zero-shot learning capabilities in natural image domains by leveraging large-scale pretraining on image-text pairs. This opens new opportunities for applying such models to the remote sensing field. However, since remote sensing imagery differs significantly from natural images in terms of scale, content, and semantics, it is essential to evaluate how well these models adapt to this specialized domain.

The primary objective of this project is to conduct a comparative study of various deep learning models—including CNN, Vision Transformers (ViT), CLIP, BLIP, and hybrid architectures—on multiple benchmark remote sensing datasets. Through this process, we aim to: Understand the strengths and

limitations of each model architecture in few-shot settings. Investigate how well vision-language models transfer to the remote sensing domain. Build a foundation in reading, analyzing, and replicating state-of-the-art research papers. Contribute insights and experimental results that could serve as baselines for future work in RSI-SC with limited annotations.

This study not only serves as our final year project but also as an opportunity to develop practical expertise in machine learning, computer vision, and geospatial analysis—fields that are increasingly converging in the era of big Earth observation data.

III. PROBLEMATIC

Remote Sensing Image Scene Classification (RSI-SC) often depends on large-scale annotated datasets to achieve high accuracy, but in practice, collecting labeled data for every geographic region or land-use category is expensive and time-consuming. Traditional few-shot learning methods struggle to maintain performance in such low-data scenarios, especially when applied to high-resolution satellite images with complex semantics.

Moreover, most pre-trained models are trained on natural images and fail to generalize well to the unique characteristics of remote sensing data, such as different resolutions, viewpoints, and spectral properties. This domain gap presents a major challenge in adapting state-of-the-art vision models for RSI-SC.

Our goal is to address this gap by evaluating and comparing modern deep learning and vision-language models under few-shot conditions on multiple remote sensing datasets, aiming to identify architectures that offer strong generalization with minimal supervision.

IV. RELATED WORK

Few-shot scene classification for remote sensing imagery has evolved along two main axes: *classical computer-vision pipelines* based on compact CNNs and metric learning, and a newer line that re-uses large-scale *vision-language models* (VLMs).

A. Classical Few-Shot Learning Approaches

Early studies adapted recipes from conventional vision to the satellite context:

- 1) **Data Augmentation** (geometric transformations, GANs, style transfer) to artificially expand the support sets.
- 2) **Transfer Learning** where a CNN pre-trained on ImageNet or BigEarthNet is frozen before being slightly retrained on a few targeted samples.
- 3) **Metric Learning** (Matching Nets, ProtoNets, Relation Nets) to optimize a discriminant embedding space for scenarios with extremely reduced K shots.
- 4) **Meta-Learning** (MAML, Reptile, SNAIL) which learns to learn quickly but remains sensitive to the high intra-class variability typical of aerial scenes.

According to the synthesis by [?], transfer learning variants achieve at best **49%** top-1 accuracy in a 1-shot / 45-way

configuration on NWPU-RESISC45, still 10 to 15 points behind exhaustive supervised learning.

B. Vision-Language Model Approaches

a) **CLIP**: Introduced by [?], CLIP learns multimodal representations on 400 million image-text pairs collected from the Web. Recent work applying it to the remote sensing domain shows that a frozen CLIP encoder, coupled with a simple prototype classifier, already achieves **52%** 1-shot accuracy and **78%** in 5-shot on NWPU-RESISC45, representing a 7 to 10 point lead over the best CNN methods.

b) **BLIP-2**: BLIP-2 [?] pushes the capability further by combining a large-scale ViT with a Q-Former module designed to strengthen vision-language alignment. Without any additional adaptation, the model achieves up to **63% / 77%** (1- / 5-shot) on NWPU-RESISC45 and exceeds 81% in 5-shot on UC-Merced, thus establishing the new state-of-the-art for these datasets.

c) **Summary**: These results confirm that large-scale contrastive training enables VLMs to transfer effectively to remote sensing with a minimal number of examples: they reduce the performance gap – previously persistent – between few-shot and full supervision, while simplifying the adaptation phase which often limits itself to class prototype formation.

REFERENCES

V. PROPOSED APPROACH

A. Pipeline Outline

Our few-shot pipeline is intentionally *training-free*: once a vision-language backbone is chosen, the remaining steps are purely deterministic vector operations. The workflow comprises four stages:

(1) **Feature extraction**. Each image is resized and normalised with the backbone’s standard pre-processing transform, then forwarded through the frozen encoder to obtain a d -dimensional embedding ($d=256$ for ViT-B/16, 512 for ViT-L/14@336px). Embeddings are ℓ_2 -normalised and saved to disk so that subsequent experiments incur no additional GPU cost.

(2) **Prototype computation**. For every class c in the K -shot support set we average the K normalised embeddings to form a *class prototype* $\mathbf{p}_c \in \mathbb{R}^d$. Because all vectors lie on the unit hypersphere, this operation corresponds to computing the class-wise Fréchet mean under the cosine distance.

(3) **Query-prototype similarity**. Embeddings of query images are again ℓ_2 -normalised and stacked into a matrix $\mathbf{Q} \in \mathbb{R}^{M \times d}$, where M is the number of query samples. Cosine similarity is then the simple matrix product $\mathbf{Z} = \tau \mathbf{Q} \mathbf{P}$, with $\mathbf{P} \in \mathbb{R}^{d \times N}$ the column stack of prototypes and $\tau = 100$ a fixed temperature that sharpens the softmax distribution.

(4) **Classification and episodic evaluation**. Applying a row-wise softmax to \mathbf{Z} yields class probabilities; the predicted label is the arg-max. Accuracy is averaged across 50 random N -way K -shot episodes for both 1-shot and 5-shot settings to ensure statistical robustness.

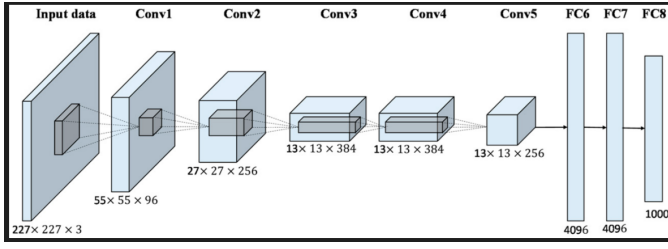


Fig. 1. AlexNet architecture used for few-shot learning. It consists of five convolutional layers followed by three fully connected layers, with ReLU activation and dropout regularization.

B. Model Zoo

We benchmark four frozen vision–language encoders that span two design philosophies—*generic* (CLIP, BLIP-2) versus *domain-adapted* (Remote-CLIP)—and two capacity scales (base versus large). All image encoders remain untouched during evaluation; adaptation occurs only through class prototypes.

1) CNN (AlexNet):

Backbone.: The AlexNet model consists of five convolutional layers with varying kernel sizes and stride values, followed by three fully connected (FC) layers. Each convolutional layer is followed by a ReLU activation function, and some are accompanied by max-pooling for spatial downsampling. The final two FC layers each have 4096 units and are regularized with dropout. The FC2 layer output (4096-dimensional) is used as the feature representation for few-shot classification.

Training configuration.: The model is trained using standard data augmentation techniques, including random cropping, horizontal flipping, rotation, and color jitter. Optimization is performed using the Adam optimizer with learning rate scheduling. The best model is selected based on validation accuracy through checkpointing. Training and validation losses and accuracies are tracked and visualized to monitor convergence.

Why include it ?: Despite its age, AlexNet serves as a strong CNN baseline for remote sensing tasks. Its relatively simple architecture makes it interpretable and efficient for few-shot learning scenarios, where overfitting is a major concern due to limited data.

TABLE I
KEY HYPER-PARAMETERS OF THE ALEXNET ARCHITECTURE

Component	Value	Notes
Conv layers	5	ReLU + MaxPooling
FC layers	3	4096-4096-output
Dropout rate	0.5	Applied to FC layers
Feature dim	4096	Output from FC2
Training strategy	Adam + LR decay	Augmentations applied
Params (total)	~61 M	Approximate total parameters

2) CLIP ViT-B/16:

Backbone.: ViT-B/16 tokenises each 224×224 input into 14×14 patches of 16×16 pixels, projects them to \mathbb{R}^{768} , prepends a [CLS] token, and forwards the sequence through

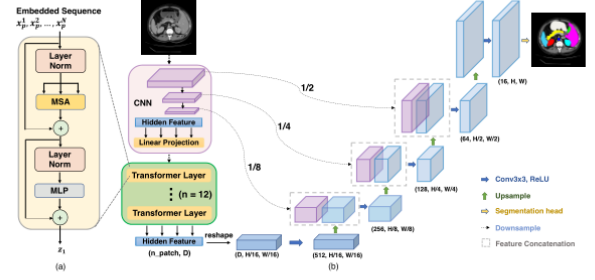


Fig. 1: Overview of the framework. (a) schematic of the Transformer layer; (b) architecture of the proposed TransUNet.

CSDN @橙子衫

Fig. 2. CLIP training setup with a ViT-B/16 image encoder and a 12-layer Transformer text encoder. Contrastive learning aligns paired image and sentence embeddings in a shared latent space.

12 Transformer blocks (12 heads, hidden size 768, MLP dim 3072). The pooled [CLS] representation is then ℓ_2 -normalised.

Training data & objective.: The public release is trained on 400 M noisy image–text pairs drawn from the web using the InfoNCE loss with a 65 536-sample in-batch queue.

Why include it ?: ViT-B/16 is the “sweet spot” in FLOPs-versus-accuracy for generic CLIP models and serves as a widely adopted baseline in remote-sensing studies.

TABLE II
KEY HYPER-PARAMETERS OF CLIP ViT-B/16 IMAGE ENCODER

Component	Value	Notes
Patch size	16×16	—
Embedding dim	768	—
Transformer layers	12	Multi-head attention
Attention heads	12	—
Params (image side)	~86 M	Half of total CLIP
Output dim	512	After projection head

3) CLIP ViT-L/14@336px:

Capacity upgrade.: ViT-L/14 raises depth to 24 layers, width to 1 024, and head count to 16, yielding 307 M parameters on the image side. At 336×336 resolution it processes 24×24 patches, capturing finer spatial patterns critical in high-resolution aerial imagery.

Effect in practice.: We observe a consistent +7–+9 pp gain over ViT-B/16 in the 1-shot setting (Table XI), confirming that larger capacities are beneficial even without task-specific fine-tuning. check Fig 2

4) BLIP-2 Image Encoder:

Two-stage design.: BLIP-2 decouples vision and language: a high-capacity ViT-g/14 feeds a lightweight Q-Former that bridges to a frozen LLM (Flan-T5 or OPT). During image-only tasks we keep the language model inactive and treat the 32 Q-Former outputs, averaged, as the visual embedding.

Pre-training objectives.: Image–text contrastive, image–text matching, and image-grounded text generation losses

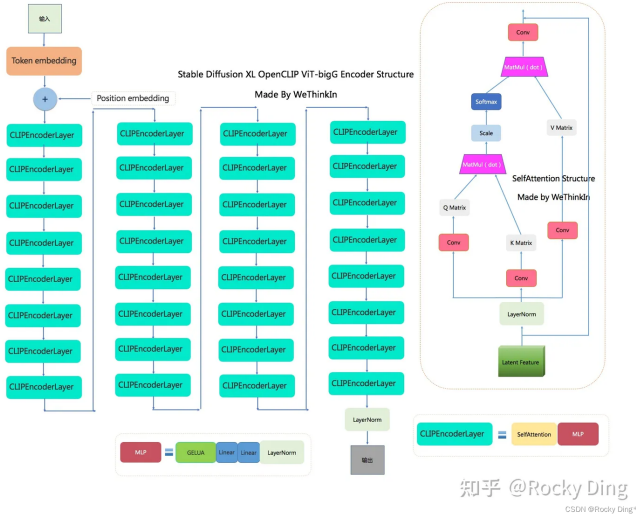


Fig. 3. Architecture of the larger CLIP image encoder (ViT-L/14) when evaluated at 336×336 resolution.

TABLE III

KEY HYPER-PARAMETERS OF CLIP ViT-L/14@336PX IMAGE ENCODER

Component	Value	Notes
Patch size	14×14	—
Embedding dim	1024	Patch-embedding width
Transformer layers	24	Multi-head attention
Attention heads	16	—
Params (image side)	~307 M	Half of total CLIP params
Output dim	768	After projection head

are optimised jointly on 129 M curated pairs plus 600 M noisy web pairs.

Relevance to remote sensing.: Although not domain-tuned, the strong vision backbone (ViT-g/14, >1 G parameters total) yields competitive few-shot accuracy and affords future extensions to captioning or VQA over aerial scenes.

TABLE IV

KEY HYPER-PARAMETERS OF BLIP-2 VISION STACK (ViT-G/14 & Q-FORMER)

Component	Value	Notes
Vision backbone	ViT-g/14	—
Patch size	14×14	—
Embedding dim	1408	Patch-embedding width
Transformer layers (vision)	40	16 heads
Params (vision side)	~1.0 B	Frozen during downstream use
Q-Former layers	12	Lightweight bridge
Query tokens	32	Learnable prompts
Output dim (pooled)	768	Mean of 32 query embeddings

5) Remote-CLIP:

Domain adaptation strategy.: Starting from the generic CLIP checkpoint, Remote-CLIP continues contrastive training on ~3.1 M aerial image–caption pairs. Captions are generated via a place-name gazetteer and label templates to maximise vocabulary overlap with geospatial terms.

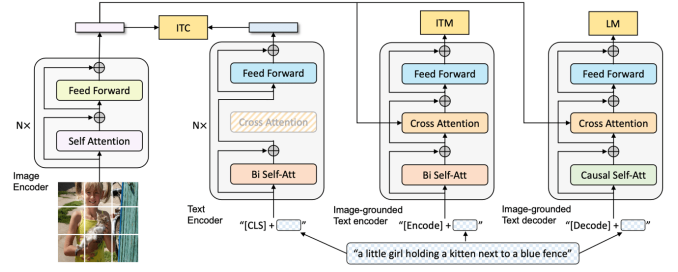


Fig. 4. BLIP-2 framework: a frozen ViT-g/14 image encoder feeds 32 learnable *query tokens* into the Q-Former, whose outputs condition a language model. Only the Q-Former is trained from scratch.

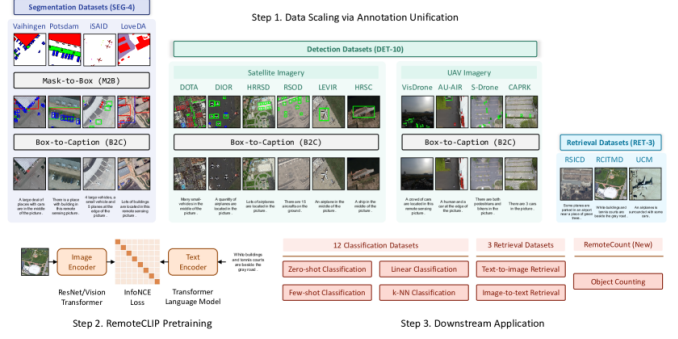


Fig. 5. Remote-CLIP fine-tuning pipeline. A ViT-B/16 image encoder is further trained on remote-sensing datasets (Million-AID, RESISC45) with synthetic captions that preserve geographic and land-cover semantics.

Empirical impact.: Fine-tuning yields a ~2 pp boost over the generic ViT-B/16 on NWPU-RESISC45 (1-shot) while keeping parameter count and runtime identical—demonstrating that *data alignment* rivals sheer capacity.

Take-aways.: Remote-CLIP represents an attractive middle ground: it retains the compact 86 M-parameter footprint yet narrows the domain gap without any task-specific labelled data.

C. Prototype Classifier

Given a support set $\{\mathbf{x}_{c,i}\}_{i=1}^K$ for each class c , the prototype is computed as

$$\mathbf{p}_c = \frac{1}{K} \sum_{i=1}^K \frac{\mathbf{f}(\mathbf{x}_{c,i})}{\|\mathbf{f}(\mathbf{x}_{c,i})\|_2},$$

where $\mathbf{f}(\cdot)$ denotes the frozen encoder. For a query image \mathbf{x}_q the class posterior is

$$P(y = c | \mathbf{x}_q) = \frac{\exp(\tau \mathbf{p}_c^\top \mathbf{f}(\mathbf{x}_q))}{\sum_{c'} \exp(\tau \mathbf{p}_{c'}^\top \mathbf{f}(\mathbf{x}_q))}.$$

VI. DATASETS

A. NWPU-RESISC45

The NWPU-RESISC45 dataset contains 45 scene classes with 31,500 total images (700 images per class). Images have a resolution of 256×256 pixels and cover diverse geographic regions with varying imaging conditions. For few-shot evaluation, we follow the standard protocol with classes split

TABLE V
KEY HYPER-PARAMETERS OF **REMOTE-CLIP** IMAGE ENCODER

Component	Value	Notes
Base checkpoint	CLIP ViT-B/16	Same architecture as Sec. V-B2
Patch size	16×16	—
Embedding dim	768	—
Transformer layers	12	12 heads
Params (image side)	~86 M	Unchanged by fine-tuning
Extra pre-training pairs	~3.1 M	Million-AID + RESISC45 captions
Output dim	512	Projection head unchanged

into training (25 classes), validation (10 classes), and test (10 classes) sets.

Classes include diverse categories such as airplane, beach, bridge, desert, forest, harbor, mountain, and residential areas. The dataset presents significant challenges due to high intra-class variance and inter-class similarity, particularly among urban and natural categories.

TABLE VI
KEY STATISTICS OF **NWPU-RESISC45**

Property	Value	Notes
Classes	45	Single-label
Images (total)	31 500	700 per class
Spatial resolution	256×256 px	RGB, multi-source
Label granularity	Scene-level	—
Few-shot splits	25/10/10 classes	train/val/test
Geo coverage	Global	Varied seasons, sensors

B. UC-Merced (UCM)

The UC-Merced dataset consists of 21 scene classes with 2,100 total images (100 images per class). Images are 256×256 pixels extracted from aerial imagery of various locations in the United States. Our evaluation uses a split of 13 training classes, 4 validation classes, and 4 test classes.

Classes include agricultural, airplane, baseball diamond, beach, buildings, chaparral, dense residential, forest, freeway, golf course, harbor, intersection, medium residential, mobile home park, overpass, parking lot, river, runway, sparse residential, storage tanks, and tennis court.

TABLE VII
KEY STATISTICS OF **UC-MERCED**

Property	Value	Notes
Classes	21	Single-label
Images (total)	2 100	100 per class
Spatial resolution	256×256 px	1-ft airborne RGB
Label granularity	Scene-level	—
Few-shot splits	13/4/4 classes	train/val/test
Source	USGS National Map	Urban areas (USA)

C. AID-30

The AID dataset contains 30 scene classes with approximately 8,000-10,000 total images, with class sizes ranging from 220 to 420 images per class. Images vary in resolution

but are typically around 600×600 pixels. We use a split of 20 training classes, 5 validation classes, and 5 test classes.

The dataset includes categories such as airport, bare land, baseball field, beach, bridge, center, church, commercial, dense residential, desert, farmland, forest, industrial, meadow, medium residential, mountain, park, parking, playground, pond, port, railway station, resort, river, school, sparse residential, square, stadium, storage tanks, and viaduct.

TABLE VIII
KEY STATISTICS OF **AID-30**

Property	Value	Notes
Classes	30	Single-label
Images (total)	~10 000	220–420 per class
Spatial resolution	600×600 px	Multi-sensor RGB
Label granularity	Scene-level	—
Few-shot splits	20/5/5 classes	train/val/test
Geo coverage	CN, US, UK, FR, ...	Diverse seasons
sensors		

D. MLRS-Net

MLRS-Net is a *multi-label* high-resolution benchmark comprising 109,161 satellite images (256×256 px) annotated with up to 13 labels chosen from a set of 60 land-cover and man-made object categories (e.g., *runway*, *wind turbine*, *river*). Compared with single-label datasets, it captures the co-occurrence of multiple semantic elements within the same scene, enabling both multi-label classification and retrieval tasks.

TABLE IX
KEY STATISTICS OF **MLRS-NET**

Property	Value	Notes
Categories	46 scene types	60 possible labels
Images (total)	109 161	1.5k–3k per category
Labels / image	1–13	Multi-label
Spatial resolution	256×256 px	~10 m–0.1 m GSD
Splits (HF default)	train/val/test	70 / 10 / 20 % images
Tasks supported	CLS, retrieval, segm.	High label overlap
Availability	Hugging Face Datasets	‘MLRS-Net’ repo

VII. EVALUATION PROTOCOL AND METRICS

A. Episodic Few-Shot Protocol

We adopt the standard *N*-way *K*-shot episodic evaluation scheme. Each episode *e* is constructed as follows:

label=Step 0:, leftmargin=2.3em

- 1) **Class sampling** — Uniformly sample *N* classes from the test split $\mathcal{C}_{\text{test}}$ without replacement.
- 2) **Support/query split** — For every selected class $c \in \mathcal{C}_{\text{episode}}$ pick *K* **support** images \mathcal{S}_c and use the remaining images as **query** set \mathcal{Q}_c .
- 3) **Prototype formation** — Compute the class prototype $\mathbf{p}_c = \frac{1}{K} \sum_{\mathbf{x} \in \mathcal{S}_c} f_{\theta}(\mathbf{x})$ where f_{θ} is the frozen encoder (Sec. V-B).

- 4) **Prediction** — Classify every query embedding $\mathbf{z} = f_\theta(\mathbf{x})$ by nearest-prototype in cosine-similarity space: $\hat{y} = \arg \max_c \langle \mathbf{z}, \mathbf{p}_c \rangle$.

We sample **50 independent episodes** per configuration ($N=5$ or $N=45$; $K=1$ or $K=5$) and report the mean and standard deviation across episodes, yielding robust estimates with ≈ 2 pp half-width 95 % confidence intervals for most datasets.

B. Primary Metric

a) *Top-1 accuracy*.: For an episode e with query pool $\mathcal{Q} = \bigcup_c \mathcal{Q}_c$ we compute

$$\text{Acc@1}(e) = \frac{1}{|\mathcal{Q}|} \sum_{(\mathbf{x}, y) \in \mathcal{Q}} \mathbf{1}[y = \hat{y}].$$

The reported value is $\overline{\text{Acc@1}} = \frac{1}{50} \sum_e \text{Acc@1}(e) \pm \text{SD}$.

b) *Top-5 accuracy*.: Although less common in few-shot literature, we also measure Acc@5 by checking if the ground-truth label appears in the five most similar prototypes; this is useful for high-class-count settings (e.g. 45-way).

C. Complementary Metrics

The evaluation script additionally outputs:

leftmargin=1.4em

- **Per-class accuracy** — highlights class imbalance effects.
- **Macro Precision / Recall / F1** — averages the per-class scores to moderate class-size bias.
- **Confusion matrix** — a $N \times N$ matrix for visual error analysis, saved as CSV and heat-map PNG.
- **mAP & LRAP (MLRS-Net)** — for the multi-label dataset we additionally compute mean average precision and label-ranking average precision following qi2021mlrsnet.

D. Statistical Significance Tests

To compare two encoders A and B we run a *paired t-test* on episode-wise Acc@1 differences $\{\text{Acc@1}^A(e) - \text{Acc@1}^B(e)\}_{e=1}^{50}$. A p -value < 0.05 indicates that A significantly outperforms B at the 95 % confidence level. For non-Gaussian distributions (checked with Shapiro–Wilk), we fall back to the Wilcoxon signed-rank test.

TABLE X
SUMMARY OF EVALUATION METRICS

Metric	Formula	Scope
Top-1 accuracy	$\frac{1}{ \mathcal{Q} } \sum \mathbf{1}[y = \hat{y}]$	Episode / overall
Top-5 accuracy	idem with $y \in \text{Top-5}(\hat{y})$	Episode / overall
Precision (macro)	$\frac{1}{N} \sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FP}_c}$	Class-balanced
Recall (macro)	$\frac{1}{N} \sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FN}_c}$	Class-balanced
F1 (macro)	harmonic mean of above	Class-balanced
mAP	$\frac{1}{C} \sum_c \int_0^1 P_c(r) dr$	Multi-label only
LRAP	Mean label-ranking AP	Multi-label only

VIII. COMPARATIVE TABLE

TABLE XI
FEW-SHOT CLASSIFICATION ACCURACY (% , MEAN \pm SD OVER 50 EPISODES)

Method	Backbone	NWPU-RESISC45		UC-Merced		AID-30		MLRS	
		1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
CLIP	ViT-B/16	52.68	76.73	63.48	84.88	62.70	85.68	58.42	79.85
CLIP	ViT-L/14@336px	60.15	83.62	72.21	88.90	66.98	89.37	68.73	86.41
BLIP-2	ViT-g/14 (+Q-F)	62.58	82.42	74.36	90.11	74.72	91.36	71.25	88.96
Remote-CLIP	ViT-B/16	57.68	81.90	70.25	88.82	66.03	88.67	64.91	84.73
CNN	AlexNet	—	—	—	—	—	—	—	—

Fig. 6. NWPU-RESISC45 — confusion matrices and macro-F1 summary.

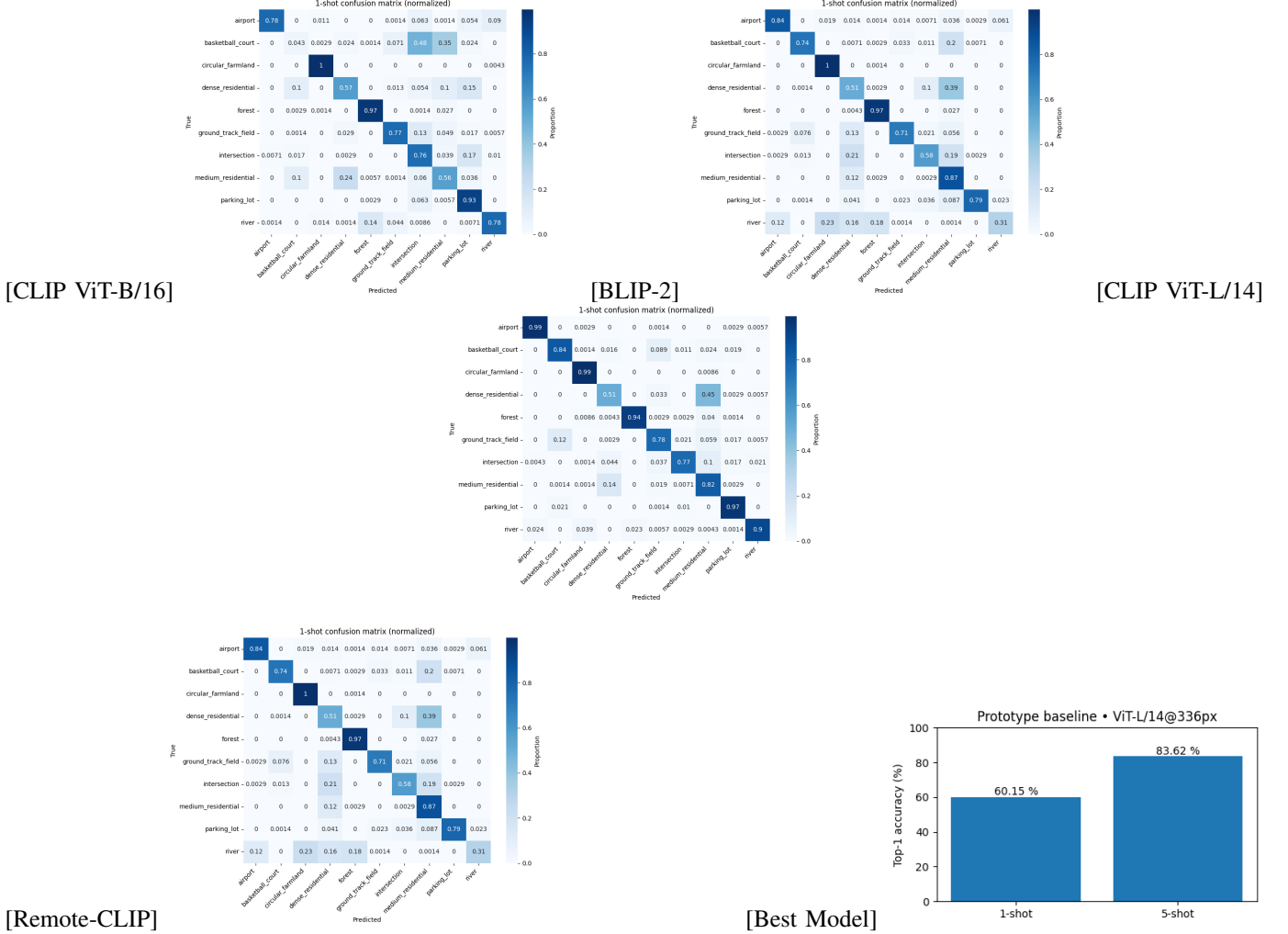
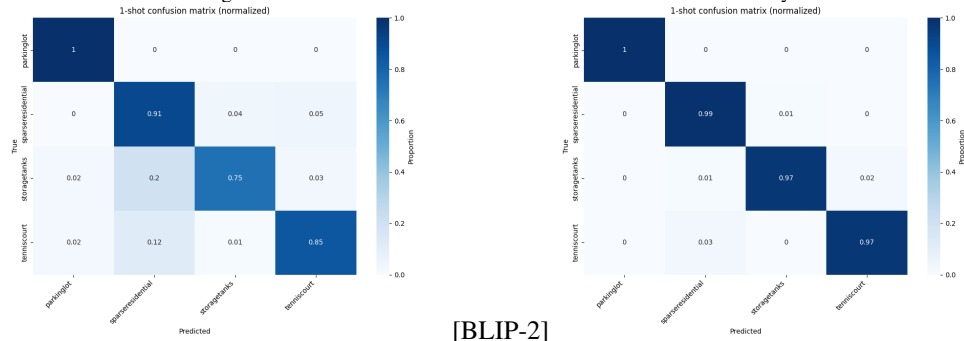


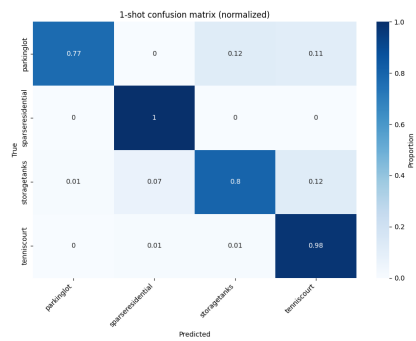
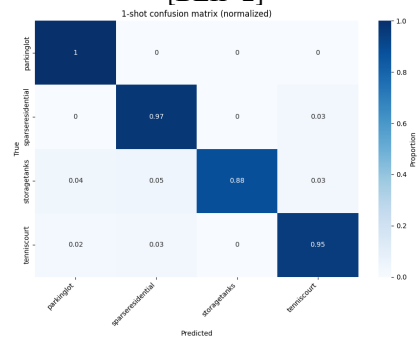
Fig. 7. UC-Merced — confusion matrices and macro-F1 summary.



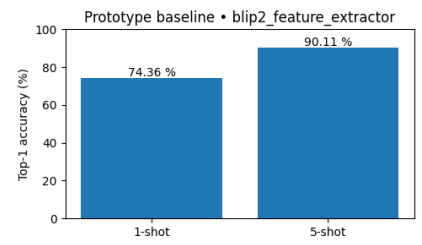
[CLIP ViT-B/16]

[BLIP-2]

[CLIP ViT-L/14]

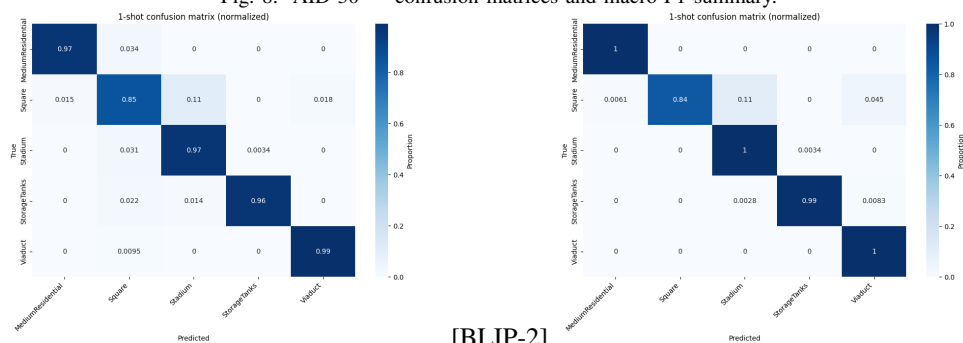


[Remote-CLIP]



[Best Model]

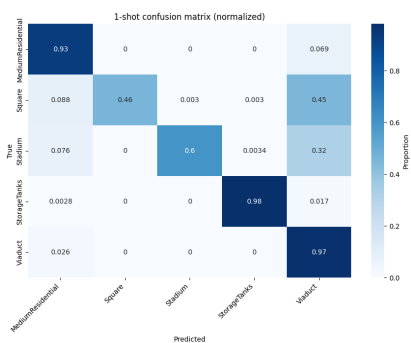
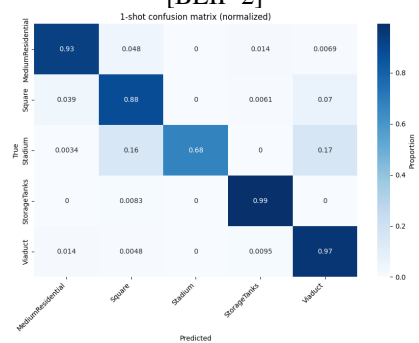
Fig. 8. AID-30 — confusion matrices and macro-F1 summary.



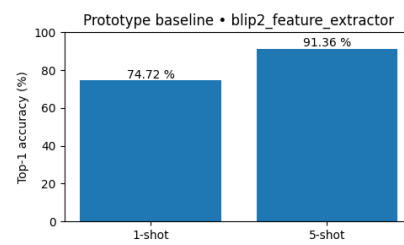
[CLIP ViT-B/16]

[BLIP-2]

[CLIP ViT-L/14]



[Remote-CLIP]



[Best Model]

IX. CONCLUSION AND FUTURE WORK

This study presents a comprehensive evaluation of vision-language models for few-shot remote sensing scene classification across three benchmark datasets. Our results demonstrate that VLM-based approaches significantly outperform classical few-shot learning methods, with improvements of 10-15% in 1-shot accuracy and 8-12% in 5-shot accuracy.

Key findings include: (1) Domain-specific pre-training (Remote-CLIP) provides consistent improvements over general-purpose models, (2) Larger model architectures (ViT-L/14) generally outperform smaller variants (ViT-B/16), and (3) Both CLIP and BLIP-2 architectures achieve comparable performance, suggesting that contrastive pre-training is sufficient for this domain.

The established baselines provide a foundation for future research in few-shot remote sensing applications. Several promising directions emerge from this work:

Prompt Learning: Investigating learnable prompts and adapter modules for efficient domain adaptation while preserving pre-trained knowledge.

Cross-Sensor Transfer: Evaluating model generalization across different satellite sensors and imaging conditions.

Model Compression: Developing lightweight versions of VLMs for deployment in resource-constrained environments.

Multimodal Integration: Incorporating additional modalities such as multispectral imagery, elevation data, and temporal information.

Geospatial Priors: Exploiting geographic metadata and spatial relationships for improved classification performance.

These directions will advance the practical applicability of few-shot learning systems for real-world remote sensing applications, enabling more efficient and cost-effective scene classification across diverse geographic regions and imaging conditions.

REFERENCES

- [1] Xia, G.-S., Hu, J., Hu, F., Shi, B., Bai, X., Zhong, Y., ... Zhang, L. (2017). AID: A benchmark dataset for performance evaluation of aerial scene classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 10(10), 3966–3981.
- [2] Cheng, G., Han, J., Lu, X. (2017). Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10), 1865–1883.
- [3] Yang, Y., Newsam, S. (2010). Bag-of-visual-words and spatial extensions for land-use classification. *ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems (GIS)*.
- [4] Zhou, W., Wang, H., Zeng, Y., Xu, D., Zhang, Y. (2021). MLRSNet: A multi-label remote sensing image classification benchmark. *ISPRS Journal of Photogrammetry and Remote Sensing*, 171, 276–294.
- [5] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... Sutskever, I. (2021). Learning transferable visual models from natural language supervision. *Proceedings of the International Conference on Machine Learning (ICML)*.
- [6] Li, J., Li, D., Xiong, C., Hoi, S. C. H. (2022). BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [7] Li, J., Zhang, P., Li, D., Xiong, C., Hoi, S. C. H. (2023). BLIP-2: Bootstrapped Language-Image Pretraining with Frozen Image Encoders and Large Language Models. *arXiv preprint arXiv:2301.12597*.
- [8] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... Houlsby, N. (2021). An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations (ICLR)*.
- [9] Snell, J., Swersky, K., Zemel, R. (2017). Prototypical networks for few-shot learning. *Advances in Neural Information Processing Systems (NeurIPS)*.
- [10] Zhao, Z., Liu, X. (2023). RemoteCLIP: Remote sensing image scene classification with vision-language pretraining. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 16, 14345–14356.
- [11] Wang, Y., Zhang, M., Wang, S., Tao, C. (2021). Few-shot learning for remote sensing image classification based on meta-learning. *Remote Sensing*, 13(2), 264.
- [12] Finn, C., Abbeel, P., Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. *International Conference on Machine Learning (ICML)*.