

DRISHTI : Deep Remote-Sensing Intelligence for Semantic Hybrid Text-Image Understanding

Abstract—The proliferation of satellite constellations and high-resolution aerial platforms has generated unprecedented volumes of Remote Sensing (RS) imagery, yet effective natural language interaction with such data remains a significant challenge. We introduce DRISHTI (Deep Remote-sensing Intelligence for Semantic Hybrid Text-Image understanding), a unified Vision-Language Model (VLM) framework enabling intuitive natural language interaction with RS imagery across varied resolutions, sensor modalities, and downstream tasks. DRISHTI addresses three critical gaps in existing RS-VLM research. First, we contribute DRISHTI-GCV, a large-scale, difficulty-aware dataset suite comprising approximately 180,000 curated samples spanning Grounding, Captioning, and Visual Question Answering (VQA) tasks, aggregated from multiple RS benchmarks with explicit annotations for image resolution, query complexity, and object scale. Second, we propose a two-stage Curriculum Learning paradigm: the first stage adapts a Low-Rank Adaptation (LoRA)-tuned VLM backbone to general RS semantics; the second stage specializes the model for fine-grained tasks using Direct Preference Optimization (DPO) alignment. Third, we introduce a dedicated numeric question handling pipeline wherein a lightweight question-type router dispatches counting and area estimation queries to a Segment Anything Model 3 (SAM3)-based module with pyramidal tiling. Comprehensive experiments on VRSBench, RSVQA, FIT-RSFG, and real-world satellite/aerial RGB, SAR, IR, and False Color Composite imagery demonstrate that DRISHTI consistently matches or surpasses state-of-the-art RS-VLMs, achieving substantial gains in caption quality (+40% BERT-BLEU), VQA accuracy (+21% over GeoChat), and counting performance while maintaining interactive latency.

I. INTRODUCTION

The advent of large-scale Earth Observation (EO) programs, including the European Space Agency’s Copernicus Sentinel constellation, NASA’s Landsat Continuity Mission, and commercial very-high-resolution (VHR) platforms such as Maxar WorldView and Planet’s SkySat fleet, has resulted in the continuous generation of petabytes of remote sensing (RS) imagery every year. This deluge of geospatial data underpins critical applications ranging from urban planning and infrastructure monitoring to disaster response, agricultural assessment, and climate change analysis. However, the traditional paradigm of manual image interpretation by domain experts cannot scale to meet increasing operational demands, motivating the development of automated and semi-automated analysis pipelines.

While sensing capabilities continue to advance rapidly, accessibility to information extracted from RS data remains limited by the need for specialized tools and technical expertise. As a result, non-expert users, including policymakers, planners, and first responders, often face delays in deriving actionable insights from satellite imagery. This growing gap

between data collection and knowledge extraction highlights the need for more intuitive and scalable mechanisms for interacting with complex geospatial data.

Recent progress in Vision-Language Models (VLMs), liu2024visual has demonstrated the feasibility of querying visual content using natural language, enabling tasks such as image captioning, visual grounding, and Visual Question Answering (VQA). Large-scale architectures such as CLIP, BLIP-2, and LLaVA (Large Language and Vision Assistant) have achieved strong performance on natural-image benchmarks by aligning textual and visual representations during pre-training. However, applying these models directly to RS imagery introduces significant challenges. Satellite images differ fundamentally from everyday photographs in scale, sensor modality, scene geometry, and object density, leading to frequent hallucination, imprecise localization, and weak numerical reasoning when existing VLMs are deployed in this domain.

Remote sensing scenes introduce several unique challenges. Data spans multiple spatial resolutions, from sub-meter optical imagery to coarse multispectral products, requiring explicit multi-scale reasoning. Scenes often contain thousands of arbitrarily oriented, densely packed objects that strain object-centric representations. Furthermore, RS imagery extends beyond RGB data and includes modalities such as Synthetic Aperture Radar (SAR), thermal infrared, and multispectral or hyperspectral sensors, each with distinct physical characteristics. Critically, many RS applications require quantitative spatial reasoning, including object counting, size estimation, and spatial relationship analysis, all of which remain underdeveloped in most off-the-shelf VLMs.

To address these limitations, we introduce **DRISHTI** (Deep Remote-sensing Intelligence for Semantic Hybrid Text-Image Understanding), a unified framework designed to bring natural language reasoning into high-resolution, multimodal satellite imagery. DRISHTI integrates language grounding, dense object understanding, and quantitative analysis within a single architecture. Rather than treating remote sensing as a simple domain adaptation problem, our approach explicitly models scale variation, object density, and geometric structure. As a result, DRISHTI enables accurate and interpretable interaction with geospatial data through natural language.

Our key contributions are summarized as follows:

- **DRISHTI-GCV: A Difficulty-Aware RS Dataset.** We introduce DRISHTI-GCV, a unified dataset of $\sim 170K$ samples for grounding, captioning, and VQA, integrating multiple RS benchmarks with annotations for resolution

(224^2 - 2048^2), query difficulty, object size, and spatial complexity. This enables curriculum learning and fine-grained evaluation across diverse RS conditions.

- **Curriculum Learning with LoRA and DPO Alignment**

We propose a two-stage training strategy where LoRA-based fine-tuning adapts the VLM to RS semantics, followed by DPO-based specialization on difficult samples using a BERT-BLEU reward to reduce hallucinations. This allows an 8B model to outperform larger baselines.

- **Grounding via the Adaptive Hierarchical Grounding Net-work (AHG-Net)**

We introduce a query-routing mechanism that directs numeric questions to a SAM3-based segmentation pipeline with multi-scale tiling and geometric post-processing, enabling accurate counting and area estimation beyond LLM-only reasoning.

II. LITERATURE REVIEW

A. Vision-Language Models for Remote Sensing

Adapting Vision-Language Models (VLMs) to remote sensing (RS) is challenging due to ultra-high resolution imagery, dense object layouts, oriented targets, and complex spatial relationships. Early RS-specific models extend general-purpose VLMs with spatial awareness and high-resolution encoding. GeoChat [1] builds upon LLaVA with a high-resolution CLIP encoder and region-based conversational grounding. LHRSGBot [2] improves spatial reasoning using curriculum-based alignment over the LHRSG-Align corpus, while SkyEyeGPT [3] unifies captioning, VQA, and grounding by injecting topic tokens into an EVA-G + LLaMA-2 backbone.

These models rely on RS-targeted datasets such as VRSBench [4], which provides image patches annotated with captions, oriented bounding boxes, and over 120k VQA pairs, complemented by RSICD [5] and RSVQA [6] for multi-resolution supervision. To address large image sizes, pyramidal tiling strategies preserve both fine detail and global context more effectively than uniform splitting, which tends to break long-range spatial dependencies.

B. High-Resolution Multimodal Representation Learning

The Qwen family provides a strong foundation for high-resolution multimodal reasoning. Qwen3-VL [7] introduces dynamic resolution handling and Multimodal Rotary Position Embeddings (M-RoPE), enabling fine-grained spatial encoding across image and video inputs. This allows the model to preserve geometric structure and object locality, which is essential in RS imagery.

Dynamic resolution enables computation to focus on relevant image regions without fully upsampling large images, improving efficiency while maintaining spatial fidelity. However, most existing VLMs still focus on RGB imagery and lack unified support for SAR, multispectral, and infrared modalities, limiting real-world RS applicability.

C. Captioning and Visual Question Answering in Remote Sensing

Remote sensing captioning and VQA demand structured semantic understanding, quantitative reasoning, and precise

spatial description. Unlike natural images, RS data involves complex terrain interactions and object distributions. Benchmarks such as RSICD [5] and VRSBench [4] facilitate training, but hallucination and spatial inconsistency remain common failure modes.

Curriculum-based instruction tuning in LHRSGBot [2] improves robustness, yet linguistic alignment remains underexplored in RS. Techniques such as Direct Preference Optimization (DPO) and BERT-BLEU scoring [8] have shown promise in improving semantic faithfulness and fluency but are rarely applied in RS-specific pipelines.

D. Visual Grounding and Open-Vocabulary Detection

Open-vocabulary grounding enables language-conditioned localization in large-scale imagery. Rex-Omni [9] introduces universal proposal networks for dense detection, while VLM-FO1 [10] connects LVLMs to detection heads. LiSAT [11] enables language-driven pixel-level reasoning for counting and area estimation, addressing key quantitative tasks in RS.

Grounding DINO [12] advances open-set detection via a cross-modal decoder and language-guided query selection, improving generalization to unseen classes. However, most grounding models struggle with cluttered scenes and rotated targets that dominate aerial imagery, motivating domain-specialized approaches.

III. DATASET PREPARATION: DRISHTI-GCV

DRISHTI-GCV was constructed to address the fragmentation of remote-sensing (RS) vision-language datasets, which typically focus on a single task, operate at fixed resolutions, or fail to model difficulty. Our dataset unifies captioning, grounding, and VQA across RGB, SAR, and infrared (IR) modalities, paired with explicit resolution metadata, difficulty annotations, and multi-scale object statistics. The resulting multi-resolution, difficulty-aware suite (Table I) supports training DRISHTI to handle tiny dense objects, heterogeneous sensor modalities, multi-resolution scenes, and numeric or spatially compositional queries.

TABLE I
DRISHTI-GCV DATASET COMPOSITION.

Task	Samples	Resolution Range
Captioning (Stage I)	~22k	256^2 - 512^2
Captioning (Stage II)	~20k	up to 2048^2
Grounding	~60k	224^2 - 2048^2
VQA (Generalized)	~30k	256^2 - 512^2
VQA (Specialized)	~20k	512^2 - 2048^2
SAR (SARLANG)	~15k	512^2
Infrared (GeoAI)	~13k	1024^2
Total	~180k	224^2-2048^2

A. Captioning Dataset Preparation

Caption annotations are sourced from VRSBench and Git10M [13]. We normalize captions through lightweight text pre-processing (tokenization, lemmatization) and filter out ultra-short or off-topic descriptions.

1) *Generalized Captioning Split (Stage I)*: The Stage I split targets broad RS semantics at moderate resolution. We select ~22k image-caption pairs (12k from VRSBench, 10k from Git-10M) spanning 256^2 - 512^2 . Qwen3-VL-32B evaluates each sample using three criteria: (i) object-caption alignment, (ii) semantic density, and (iii) n-gram repetition: yielding a difficulty score from 1 to 4. We retain samples rated 1-3 based on object size, clutter, and spatial-constraint complexity.

2) *Specialized Captioning Split (Stage II)*: Stage II focuses on complex, high-density scenes. We oversample images with many small objects, strong spatial relations, and high visual clutter. Most retained captions fall within difficulty levels 3-4, and we additionally include upsampled high-resolution tiles ($\geq 1024^2$).

B. Grounding Dataset Preparation

Our grounding suite (~30k samples) integrates VRSBench (512×512 tiles with oriented bounding boxes (OBBs) and referring expressions), OPT-RSVG (224×224 and 800×800 crops) [14], RSVG/RSVG-HR [15], and DIOR/DOTA-v2 [16], [17], spanning resolutions from 224^2 to 2048^2 .

DIOR and DOTA provide only a few object classes and lack referring expressions; we therefore generate natural-language descriptions using GPT-5-mini to unify grounding supervision across datasets. To standardize localization quality, we construct consistent OBBs where originals are absent or axis-aligned. A hybrid pipeline combining SAM-3 box-prompted segmentation with classical CV methods (contours and `minAreaRect`) yields rotation-aware bounding boxes for all samples, ensuring coherent spatial annotations suitable for multi-scale grounding.

C. VQA Dataset Preparation

For VQA, we integrate VRSBench with RSVQA-LR/HR to build a difficulty-aware, multi-resolution QA corpus. VRSBench provides rich question-answer structures that capture spatial context and reasoning depth, while RSVQA adds both low- and high-resolution variants. Difficulty scores (1-5) assigned by Qwen3-VL-32B reflect constraint count, spatial-relational complexity, and required domain knowledge.

1) *Generalized VQA Split*: The generalized split comprises ~30k QAs (9k numeric, 10.5k classification, 10.5k semantic), retaining only difficulty 1-3 samples for broad-coverage training.

2) *Specialized VQA Split*: The specialized split (~20k QAs) focuses on difficulty 4-5 questions. These samples come from high-resolution aerial tiles ($\geq 512^2$) where tiny structures remain resolvable. This subset targets fine-grained relational reasoning under clutter.

D. SAR and IR Dataset Preparation

We initially explored cross-modal translation methods such as Pix2Pix [18], ThermalGAN, and ControlNet-style diffusion [19] to convert SAR and IR/thermal imagery into RGB for reuse of natural-image encoders. However, these transformations proved unstable: they frequently hallucinated textures

or over-smoothed gradients, degrading modality-specific structure critical for RS understanding.

1) *SAR Preparation*: To preserve SAR characteristics, we directly subsample 15k captioning and 15k VQA samples from SARLANG-1M [20], converting them into instruction-format pairs for multimodal alignment without altering radiometric properties.

2) *Infrared (IR) Preparation*: For IR, we prioritize principled dataset construction. Using Microsoft Planetary Computer’s GeoAI interface, we retrieve geographically anchored scenes from the National Agriculture Imagery Program (NAIP), covering diverse U.S. regions at 0.3-1 m GSD with RGB+NIR bands. We extract all spectral channels and generate true-color, NIR-only, and false-color composites, followed by tiling into 1024×1024 crops to capture both global context and fine structures. From over 300 scenes, we assemble a large-scale IR corpus and annotate it with a GPT-5-mini prompting pipeline tailored to RS semantics, followed by selective human verification. This yields 12k VQA instances and 2k captions. Unlike cross-modal translation, this modality-specific approach preserves radiometric integrity and significantly improves captioning and VQA performance across IR modalities.

IV. TRAINING AND INFERENCE PIPELINE

We adopt a two-stage, curriculum-style pipeline built on Qwen-3-VL-8B-Instruct, chosen for its native support for high-resolution inputs via dynamic tiling and resolution-aware visual encoding, which matches the multi-scale nature of RS imagery. Stage I adapts the model to remote sensing using DRISHTI-GCV generalized mixed-task dataset (captioning + VQA), while Stage II performs task-specific specialization for captioning and VQA, including DPO-based alignment and a SAM3-powered, tool-augmented numerical Grounding module. Figure below summarizes the overall training pipeline.

A. Stage I: General Remote-Sensing Adaptation

The pretrained model **Qwen3-VL-8B-Instruct** is adopted as the base vision-language backbone, chosen for its support of high-resolution inputs and multimodal positional embeddings. A mixed remote-sensing dataset comprising image-caption pairs and image-question-answer triples from DRISHTI-GCV’s generalized split is used for supervised fine-tuning. The adaptation is performed via **LoRA**, updating only low-rank parameter matrices in attention and MLP layers, while freezing the base weights. At the end of this stage, a domain-adapted checkpoint is obtained, which serves as the foundation for all subsequent task-specific specialization.

B. Stage II-A: Caption-Specialized Curriculum

Task-specific SFT. From the generalized checkpoint, we specialize a captioning branch using the DRISHTI-GCV specialized caption split, which emphasizes semantically dense descriptions, multiple object types and relations, and multi-scale scenes (tiny objects to city-scale context) up to $\sim 2048^2$ resolution. We reuse the Stage I LoRA setup and optimizer, minimizing autoregressive cross-entropy over caption tokens.

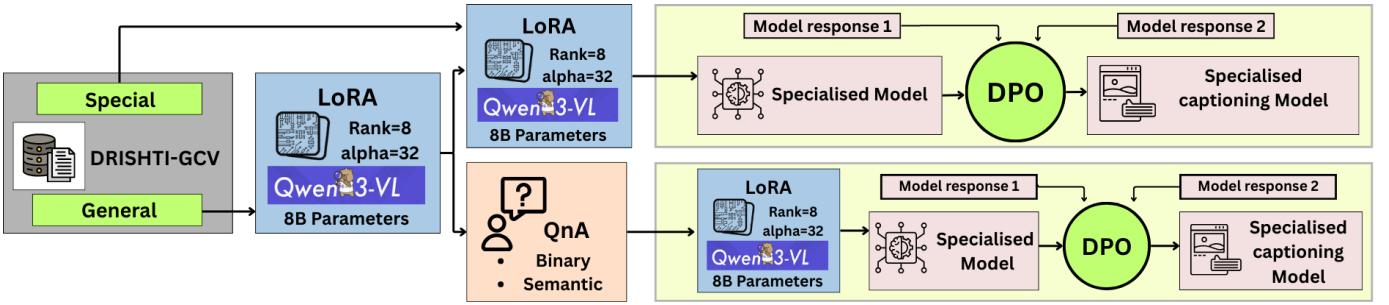


Fig. 1. Training Pipeline

C. Stage II-B: VQA-Specialized Curriculum

The VQA module proceeds in three steps:

- 1) **Specialized fine-tuning.** The RS-adapted backbone is fine-tuned on specialised subset of DRISTI-GCV covering semantic, yes/no, and numeric (counting/area) questions.
- 2) **Query routing.** At inference, a lightweight classifier (Qwen3-VL-30B) routes each question into one of three branches: *Numeric* → segmentation-based counting/area module, *Yes/No* → VQA head, *Semantic* → VQA head with optional preference-based refinement.
- 3) **Numeric reasoning.** For numeric queries: (i) generate a short segmentation prompt from the question, (ii) apply multi-scale segmentation via a tiling module + promptable segmenter (e.g. SAM3), (iii) fuse overlapping masks, (iv) compute object count and physical area (via GSD × pixel count), (v) return numeric answer in natural language.

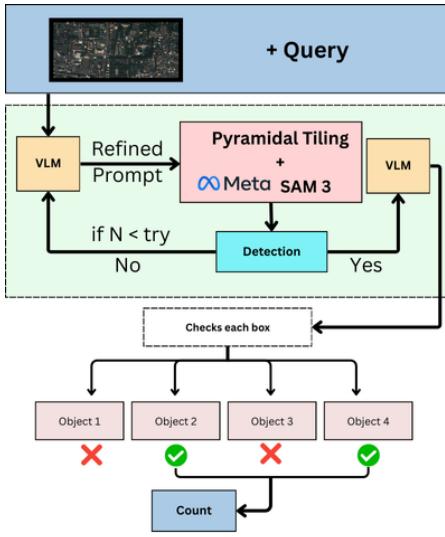


Fig. 2. Counting/Area Pipeline

D. Grounding via the Adaptive Hierarchical Grounding Network (AHG-Net)

We introduce the *Adaptive Hierarchical Grounding Network* (AHG-Net), a specialized architecture designed to bridge the

gap between open-vocabulary segmentation and nuanced spatial reasoning. Addressing the challenges of scale variation and high-density mask outputs, the AHG-Net integrates a centroid-based optimization mechanism with a stochastic refinement policy. The grounding pipeline proceeds through five distinct stages:

- 1) **Subject-Reference Query Decomposition:** Complex queries often involve relative positioning that conflates target objects with their environment. To disambiguate these roles, we first decompose the input query q into two semantic constituents:
 - **Primary Subject (ϕ_{pri}):** The specific entity to be located (e.g., “the car to the left of a bridge”).
 - **Secondary Subject (ϕ_{sec}):** The reference object establishing the spatial context (e.g., “car to the left of a bridge”).

This decomposition allows the system to treat the secondary subject purely as a geometric anchor, focusing the computationally expensive refinement solely on the primary target.

- 2) **Pyramidal SAM3 Segmentation:** Standard segmentation frequently misses small objects or fragments large ones. We address this via a pyramidal tiling strategy using the SAM3 encoder. The image is processed at multiple discrete resolution levels to generate a comprehensive candidate pool for both subjects:

$$\mathcal{M}_{\text{raw}} = \{(m_j, s_j)\}_{j=1}^J, \quad s_j \in [0, 1].$$

This ensures that the initial proposal set captures both dominant structural elements and fine-grained details.

- 3) **Centroid-based Clustering for Dense Segmentation Masks:** In high-resolution settings, the pyramidal encoder often generates an excessive number of redundant candidates ($J \gg T_{\max}$), creating a bottleneck for downstream processing. To optimize this, we employ a density-aware clustering strategy:

- **Greedy box clustering:** Candidate bounding boxes are clustered in the 2D coordinate space using a greedy, centroid-based procedure. Starting from a randomly sampled box as the initial centroid, each remaining box is assigned to the nearest existing centroid (in terms of center-distance or $1 - \text{IoU}$);

if this distance exceeds a threshold, a new cluster is created. This process yields clusters $\{C_k\}$ of geometrically similar candidates.

- **Centroid selection:** For each cluster C_k , a representative *centroid mask* m_k^* is chosen by minimizing the distance to the cluster mean:

$$m_k^* = \arg \min_{m \in C_k} d(c(m), \bar{c}_k),$$

where $c(m)$ denotes the box center of m and $\bar{c}_k = \frac{1}{|C_k|} \sum_{m \in C_k} c(m)$ is the mean center of cluster C_k . This greedy centroid selection effectively selects closer candidates within each cluster.

This step splits the large number of input proposals into smaller optimized sets of potential Subjects that can then be classified by the VLM.

- 4) **Subject Confidence Estimation:** Within each cluster, boxes are viewed in the context of the full image and are classified into three bands — *High*, *Medium*, or *Low* confidence of being query-relevant subjects. High-confidence boxes whose structure and colour attributes closely match the described subject class are immediately added to a provisional *selected mask set*, while the remaining Medium- and Low-confidence boxes are retained for further refinement.
 - 5) **Chunked VLM refinement over ambiguous proposals:** The union of Medium- and Low-confidence boxes is partitioned into manageable chunks of size K . For each chunk, the VLM is provided only the corresponding bounding-box crops and the textual query. It re-evaluates each candidate as High, Medium, or Low confidence subject. Low-confidence boxes are discarded. Medium-confidence boxes are then examined at finer scale by cropping around each box and querying the VLM again, allowing it to focus on local structure, colour, and other object attributes. Newly promoted High-confidence boxes from this stage are added to the selected mask set, progressively filtering out erroneous detections while retaining plausible subjects (both primary and secondary).
 - 6) **Final spatial reasoning over primary and secondary subjects:** After the refinement loop, the selected mask set contains only subject-consistent candidates. The VLM orchestrator then performs spatial reasoning to determine which of these masks correspond to the primary and secondary objects specified in the query. Primary objects may be defined globally (e.g., “the northernmost stadium”) or relationally with respect to secondary objects (e.g., “the ship to the left of the pier near the refinery”). Using these directional and relational cues, conflicting masks are eliminated and the subset of masks that best satisfies the query semantics is retained as the final grounded solution.
- Since evaluating all ambiguous masks via a Vision-Language Model (VLM) is computationally prohibitive, we employ a **stochastic batching strategy**:

$$\mathcal{B}_k \sim \text{Uniform}(\mathcal{M}_{\text{amb}}).$$

Random subsets \mathcal{B}_k are sampled and fed to the multimodal controller, **Qwen3-VL-30B-A3B-Instruct** for classification. The VLM functions as a semantic discriminator, analyzing the visual context, the spatial anchor \tilde{m}_{sec} , and the batch simultaneously. It identifies the subset of masks within \mathcal{B}_k that optimally align with the primary query ϕ_{pri} , effectively resolving ambiguity through contextual reasoning.

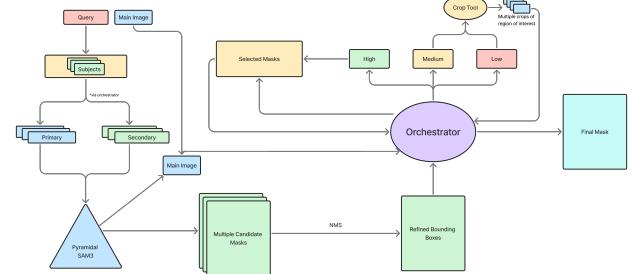


Fig. 3. SAM3 based Grounding pipeline with Qwen3-VL Orchestrator

E. FCC-to-RGB Reconstruction

We start from four-band NAIP imagery (R, G, B, NIR) and generate synthetic false-color composites (FCCs) by randomly sampling any 3-band subset that includes NIR, then permuting them.

A dual-path classifier predicts the spectral identity of each FCC channel, assigning each channel to one of {R, G, B, NIR}. Let the predicted assignments for the three channels be $\hat{y}_1, \hat{y}_2, \hat{y}_3$. The missing RGB band is then identified as:

$$B_{\text{miss}} = \{R, G, B\} \setminus \{\hat{y}_1, \hat{y}_2, \hat{y}_3\}.$$

If red or green is missing, we reconstruct it using a root-polynomial color correction (RPCC) mapping. For a pixel with input color vector $x = [x_1, x_2, x_3]^T$, the corrected channel output is:

$$C_{\text{out}} = \psi(x)^T v,$$

$$\psi(x) = [x_1, x_2, x_3, \sqrt{x_1 x_2}, \sqrt{x_1 x_3}, \sqrt{x_2 x_3}, 1]^T.$$

Here, v denotes the learned coefficient vector (from calibration). RPCC is known to better preserve color under varying exposure compared to conventional polynomial regression. If the blue channel is missing instead, we use a pretrained 3-D lookup table (LUT) with trilinear interpolation to estimate the missing band:

$$B_{\text{out}} = \sum_{i=1}^8 w_i T_i,$$

where T_i are the eight surrounding LUT vertices and w_i are the interpolation weights.

The final reconstructed RGB image is

$$I_{\text{RGB}} = \{R_{\text{out}}, G_{\text{out}}, B_{\text{out}}\},$$

which is then fed to downstream vision-language pipelines.

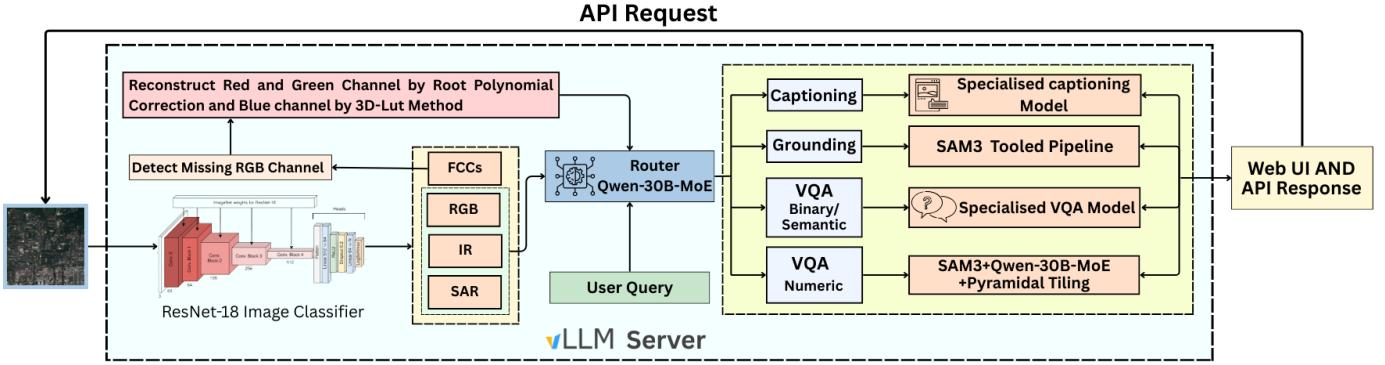


Fig. 6. DRISHTI System Architecture: A microservices-based orchestration layer dispatches multimodal queries to specialized inference engines via an intelligent routing fabric.

VII. SYSTEM DESIGN AND USER INTERFACE

A. Architectural Overview

The DRISHTI system employs a modular inference architecture designed to handle the high variability of remote sensing queries. The core pipeline (Figure 6) orchestrates interactions between a lightweight routing layer, specialized vision-language models (VLMs), and geometric segmentation tools.

Routing Logic. To optimize computational resources, we implement a dual-routing strategy.

- **Visual Routing:** A lightweight ResNet classifier categorizes input imagery (e.g., dense urban, maritime, or SAR) to pre-filter eligible downstream tools.
- **Query Routing:** A frozen Qwen3-VL-30B-A3B serves as an “LLM-as-a-judge,” analyzing user intent to dispatch requests to one of three specialized pipelines: *Captioning*, *Grounding*, or *VQA*. Within *VQA*, the router further discriminates between semantic queries (routed to the VLM) and numeric/spatial queries (routed to the segmentation engine).

Inference Engine. The system leverages **Qwen3-VL-8B** as the primary reasoning backbone, fine-tuned via the DRISHTI curriculum with LoRA adapters. To support high-precision object counting and area estimation, we integrate **SAM3** (Segment Anything Model 3). SAM3 operates in a pyramidal tiling mode, allowing it to detect small objects across large gigapixel satellite scenes. These geometric outputs (masks, oriented bounding boxes) are injected back into the VLM context for grounded reasoning.

B. Implementation and Interface

The backend is built on a scalable microservices pattern. All VLMs are served via a vLLM-based inference server, enabling optimized KV-cache reuse and batching across different tools. The frontend, implemented in React, provides a unified multimodal chat interface (Figure 8) where users can upload RGB, SAR, or IR imagery.

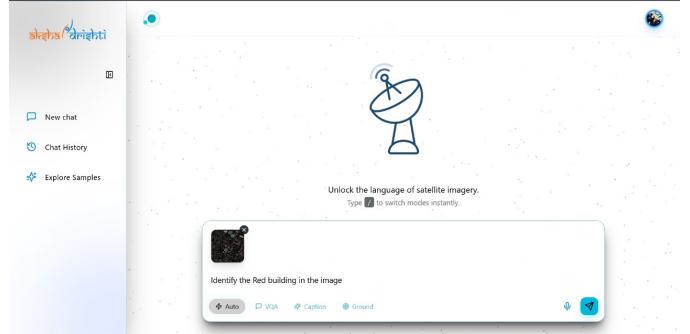


Fig. 8. The DRISHTI user interface featuring a unified chat panel for multimodal interaction and visualization of model outputs.

Crucially, the system maintains a “contextual memory” of intermediate artifacts. If a user asks to “count the ships” (triggering SAM3) and subsequently asks “are they moving?”, the system retrieves the previously computed bounding boxes and masks from cache rather than re-running segmentation, ensuring low-latency conversational interactivity.

VIII. CONCLUSION AND FUTURE WORK

This paper presented DRISHTI, a unified remote sensing vision–language framework combining a curriculum-trained Qwen3-VL backbone with SAM3-based segmentation. By decoupling numeric reasoning from semantic generation, the system achieves robust performance across diverse resolutions and modalities. Future work will focus on integrating **Complex-Valued CNNs** into the VLM encoder to natively process SAR phase and amplitude data, a novel architectural shift that unlocks physical scattering properties currently ignored by standard real-valued transformers. We also plan to extend the framework to multi-date temporal reasoning.

REFERENCES

- [1] K. Kuckreja, M. S. Danish, M. Naseer, A. Das, S. Khan, and F. S. Khan, “Geochat: Grounded large vision–language model for remote sensing,” *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [2] D. Muhtar, Z. Li, F. Gu, X. Zhang, and P. Xiao, “Lhrs-bot: Empowering remote sensing with vgi-enhanced large multimodal language model,” *arXiv preprint*, 2024.

TABLE XIII
DRISHTI-GCV SOURCE DATASET BREAKDOWN.

Source	Samples	Tasks
VRSBench	29,614	Caption, VQA, Ground
Git-10M (filtered)	10,000	Caption
RSVQA-LR	8,500	VQA
RSVQA-HR	7,200	VQA
OPT-RSVG	12,000	Grounding
RSVG/RSVG-HR	4,800	Grounding
DIOR/DOTA-v2	6,200	Grounding

TABLE XIV
OBJECT CATEGORY DISTRIBUTION IN DRISHTI-GCV.

Category	Percentage
Vehicles (cars, trucks, planes)	28.3%
Buildings (residential, commercial)	24.1%
Infrastructure (roads, bridges)	18.7%
Water bodies & ships	12.4%
Agricultural land	9.2%
Other (parks, forests, etc.)	7.3%

Through this objective, the model is aligned to follow multi-modal instructions and to produce coherent, task-appropriate outputs conditioned on both the image and the text prompt.

SYSTEM PROMPT (BRIEF VERSION)

You are a remote-sensing segmentation and mask-verification agent operating under a strict one-shot workflow. You get exactly one call to `segment_phrase_batch` and must include 3–5 true synonyms of the target feature (add color only if explicitly stated). SAM3 returns unique masks via NMS; if zero masks are produced, immediately call `report_no_mask()`.

After segmentation, classify each mask by quick-scan confidence: HIGH (clear match), MEDIUM (likely but uncertain), LOW (ambiguous or possibly wrong). HIGH confidence masks may be accepted immediately. MEDIUM or LOW confidence masks must be examined with `examine_each_mask()` before selection. Determine whether masks represent the primary target, not reference or intermediate objects; if all masks are the wrong feature type, you must output `report_no_mask()` because the one-shot chance is already spent.

Perform all spatial reasoning using mask bounding-box centroids, interpreting terms like “leftmost,” “in the left,” “below,” “near,” or compounds such as “below and rightmost of X.” Apply directional filtering first when needed, then superlative selection. Always ensure that selected masks satisfy both feature type and spatial constraints.

For dark masks, apply mandatory shadow vs. real-object discrimination. Shadows show elongation, uniform darkness, adjacency to tall objects, and consistent sun-angle direction; real objects have structure, edges, and shape consistent with the target. After confidence, target, spatial, and shadow checks, call either `select_masks_and_return`,

`examine_each_mask`, or `report_no_mask` with no extra text after the tool call.¹

% Decrements the counter if you removed it, to avoid skipping a number (optional)

¹**Query Refactoring and Grounding:** Assisted by **ChatGPT GPT-5 Mini** (OpenAI). The model was utilized for its ability to refactor initial queries into precise, grounded search terms for enhanced information retrieval. Usage is subject to OpenAI’s applicable Terms of Use and pricing policies.

Rephrasing and Clarity: Assisted by **Grammarly Students (Premium)** for advanced rephrasing and improved clarity of the final text content. Note: This tool was used on the text content itself, not for native command-level editing within the `LATEX` source files.

Queries	REX-Omni	LISAT	Ours	Ground Truth
Leftmost Storage Tank				
Largest ship in the right half				

TABLE XV

VISUAL COMPARISON OF REX-OMNI, LISAT, AND THE PROPOSED METHOD AGAINST THE GROUND TRUTH FOR THREE DIFFERENT VISUAL GROUNDING QUERIES

Query	Image	GPT-4o-mini	Qwen3-VL-32B-Instruct	Ours	Ground Truth
In which direction is the sun facing?		south west	top right	north west	north west
How long does the game go on for in the image?		60 minutes	30 minutes	48 minutes	48 minutes

TABLE XVI

COMPARISON OF RESULTS OF GPT-40-MINI,
 QWEN3-VL-32B-INSTRUCT, AND THE PROPOSED METHOD AGAINST THE
 GROUND TRUTH FOR THREE DIFFERENT VISUAL QUESTION ANSWERING
 QUERIES