SARChat-Bench-2M: A Multi-Task Vision-Language Benchmark for SAR **Image Interpretation**

Anonymous ACL submission

Abstract

As a powerful all-weather Earth observation tool, synthetic aperture radar (SAR) remote sensing enables critical military reconnaissance, maritime surveillance, and infrastructure 005 monitoring. Although Vision language models (VLMs) have made remarkable progress in natural language processing and image understanding, their applications remain limited in professional domains due to insufficient domain expertise. This paper innovatively proposes the first large-scale multimodal dialogue dataset for SAR images, named SARChat-2M, which contains approximately 2 million high-quality image-text pairs, encompasses diverse scenarios with detailed target annotations. This dataset not only supports several key tasks such as visual understanding and object detection tasks, but also serves as the first visual-language benchmark in the SAR domain. Through this work, we enable and evaluate VLMs' capabilities in SAR image interpretation, providing a paradigmatic framework for constructing multimodal datasets across various remote sensing vertical domains. Through experiments on 16 mainstream VLMs, the effectiveness of the dataset has been fully verified. The project will be released at https://anonymous.4open. science/r/SARChat-D0ED/.

Introduction 1

007

011

017

019

027

037

041

In recent years, deep neural networks, notably CNNs (LeCun et al., 1998) and ViTs (Dosovitskiy, 2020), have achieved remarkable progress in remote sensing data analysis, enhancing both processing efficiency and analytical accuracy. However, existing research mainly focuses on visual feature extraction, while lacking deep semantic parsing and reasoning capabilities (Li et al., 2024), limiting model applicability in complex scenarios.

With the advancement of Large-Language Models (LLMs), Vision-Language Models (VLMs),

through integrating pre-training and instruction tuning, have demonstrated robust zero-shot learning and generalization in multimodal tasks (Dai et al., 2023). This has inspired researchers to explore the deep integration of visual models with LLMs.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

Although models designed for optical remote sensing images, like RSGPT (Hu et al., 2023) and GeoChat (Kuckreja et al., 2024), have shown preliminary achievements, they struggle to perform well in SAR applications. SAR images inherently pose significant interpretation challenges due to their scattering imaging mechanisms, characterized by blurred target edges, dispersed speckles, and orientation sensitivity. Meanwhile, existing SAR datasets primarily focus on visual recognition tasks (Kuckreja et al., 2024; Cheng et al., 2022; Zhang et al., 2023), leaving a critical shortage of largescale, high-quality image-text alignment datasets. Both these intrinsic characteristics and data limitations impede the advancement of VLMs in the SAR domain.

Current VLMs are primarily trained on conventional natural images without extensive fine-tuning for the SAR vertical domain. Despite their strong visual capabilities for natural images, these VLMs still have significant room for improvement in SAR image interpretation. Building upon the SARDet-100K dataset (Dai et al., 2024) with its rich SAR imagery and detection annotations, we construct SARChat-Bench-2M, a task-oriented SAR-specific image-text pair dataset, to address the insufficient SAR image interpretation capabilities of existing VLMs.

As shown in Figure 1, we present SARChat-2M, a large-scale multimodal conversational dataset for SAR images, and establish SARChat-Bench, a comprehensive multimodal task-oriented benchmark for the SAR domain. The SARChat-2M dataset contains approximately 2 million high-quality SAR image-text pairs across maritime, terrestrial, and urban scenarios, featuring fine-grained semantic de-



Figure 1: An overview of SARChat-Bench-2M. The left figure demonstrates the representative tasks realized with the SAR image-text dataset, SARChat-2M, constructed in this paper. Validating the dataset's efficacy and superiority in supporting multi-task applications. The right figure presents the correlation radar charts and quantitative line graphs derived from the performance evaluation of 16 VLMs basing on this dataset, establishing the benchmark (SARChat-Bench) within this domain.

scriptions and multi-scale resolutions (0.3-10 me-083 ters). The dataset supports major vision-language tasks such as image captioning, VQA(Visual Question Answering), visual localization, and object 086 detection. To systematically evaluate model performance in these domains, we design six specific benchmark tasks in SARChat-2M: classification, fine-grained description, instance counting, spatial grounding, cross-modal identification, and referring. To validate the effectiveness of our dataset and benchmark, we conduct extensive experiments by fine-tuning 16 state-of-the-art VLMs of varying parameter scales, including InternVL2.5, DeepSeekVL, GLM-Edge-V, and the mPLUG-Owl family. Through training on SARChat-2M, these visual language models (VLMs) acquire comprehensive multi-task capabilities in SAR interpretation, as demonstrated by our systematic evaluation on 100 SARChat-Bench. 101

091

102

104

105

106

107

109

110

111

112

113

114

The primary contributions of this paper are as follows:

- 1. The construction of SARChat-2M, the largest SAR remote sensing instruction-following dataset to date, comprising over 2 million high-quality image-text pairs across multiscenario task-oriented dialogues, alleviating the knowledge scarcity of VLMs in the SAR domain.
- 2. The development of SARChat-Bench, a comprehensive SAR domain multimodal benchmark encompassing six core tasks (classification, description, counting, localization, recog-

nition, and refering), enabling systematic evaluation of vision-language models through multi-dimensional assessment metrics.

115

116

117

118

119

120

121

122

123

124

125

126

127

3. It pioneers a research paradigm applicable to the SAR field, providing reference ideas for the construction of models in other remotesensing vertical domains. The methods and processes adopted in data collection, annotation, as well as model training and evaluation in this study have good generality and extensibility.

Related Work 2

VLMs for Remote Sensing 2.1

VLMs are capable of converting images into natu-128 ral language descriptions and parsing the relation-129 ships between objects, demonstrating remarkable 130 performance in tasks such as text-image retrieval, 131 image captioning, and visual question answering. 132 Recently, models like RemoteClip (Liu et al., 2024) 133 have been applied to the field of remote sensing im-134 ages, primarily focusing on cross-modal retrieval and zero-shot classification. However, these mod-136 els have not addressed tasks such as image de-137 scription generation and visual grounding. The 138 RSGPT model has achieved text description and 139 visual question answering for remote sensing im-140 ages, but it has not expanded to tasks such as clas-141 sification and detection. The GeoChat model has 142 advanced multi-task conversational processing of 143 high-resolution remote sensing imagery, includ-144



Figure 2: **Construction of SARChat-2M dataset.** On the left, ten existing SAR detection benchmark datasets. The middle part is the SARDet-100K dataset, formed by integrating the ten datasets on the left. On the right, six core tasks constructed based on the dataset are presented, with each task corresponding to different task identifiers, operation steps, and relevant templates.

ing scene classification, visual question answering, 145 146 multi-turn dialogue, visual grounding, and reference object detection. However, these models, in-147 cluding GeoChat, predominantly rely on optical 149 remote sensing training data, leading to suboptimal performance in SAR-specific interpretation tasks. 150 EarthGPT (Zhang et al., 2024) has extended the ap-151 plication of multimodal large language models to the remote sensing field through instruction tuning, 153 but its performance in SAR image multi-task pro-154 cessing still needs improvement. Compared with 155 natural images, the interpretation of SAR images 156 is more challenging, which poses higher demands on the model's processing capabilities and adapt-158 ability. 159

2.2 Remote Sensing Vision-Language Datasets

160

Remote sensing datasets are essential for models 161 that interpret remote sensing imagery. Existing 162 datasets such as UCM Captions (Qu et al., 2016a), 163 Sydney Captions (Qu et al., 2016b), RSICD (Lu et al., 2017), RSITMD (Yuan et al., 2022), and 165 RSVG (Zhan et al., 2023) provide preliminary re-166 sources for studying the correlation between remote sensing images and text. However, these 169 datasets are limited not only in scale but also in modality, containing only optical images without SAR data, leaving SAR interpretation ca-171 pabilities largely unexplored. Although large-172 scale datasets like MillionAID (Long et al., 2021), 173

FMoW (Christie et al., 2018), and BigEarthNet (Sumbul et al., 2019) exist, they lack text-image pairs. The RS5M dataset (Zhang et al., 2023), containing 5 million image-text pairs, is still limited to optical images. The MMRS-1M dataset (Zhang et al., 2024), which covers optical, infrared, and SAR modes, has a very low proportion of SAR image-text data. Therefore, this paper constructs the SARChat-2M dataset, which focuses on SAR images and contains over 2 million image-text pairs, covering tasks such as classification, detection, caption generation, VQA, and visual grounding.

174

175

176

177

178

179

180

181

182

183

184

185

187

189

190

191

192

193

194

195

196

197

198

199

201

3 Data Construction and Description

3.1 The Procedure of Data Construction

3.1.1 Dataset Overall

As shown in Figure 2, we propose SARChat-2M, a multi-task dataset for SAR images, comprising 2 million multimodal dialogue samples (1,836,912 train and 226,636 test samples) to ensure robust model training and evaluation.

Based on the SARDet-100K dataset (Dai et al., 2024), it incorporates multimodal adaptations and enhanced language annotations from ten established SAR detection benchmarks such as AIR-SARShip(1.0&2.0) (Wang et al., 2019a), HRSID (Wei et al., 2020), MSAR (Chen et al., 2022), SADD (Zhang et al., 2021a), SAR-AIRcraf (Zhirui

303

253

254

et al., 2023), ShipDataset (Wang et al., 2019b), SSDD (Zhang et al., 2021b), OGSOD (Wang et al., 203 2023), and SIVED (Lin et al., 2023). The SARChat-204 2M covers six semantic categories (ships, tanks, bridges, ports, aircraft, and automobiles) and supports six core SAR image analysis tasks: classifica-207 tion, fine-grained description, instance counting, spatial grounding, cross-modal identification, and referring. These diverse tasks are designed to 210 enhance VLMs' capabilities in SAR image interpre-211 tation, with 2 million carefully curated annotations for cross-modal learning. 213

3.1.2 Task Definition

214

215Based on the characteristics of SAR images and the216core capabilities of the VLM, this study constructs217an evaluation system consisting of six tasks. The218definitions of each task are as follows:

(1) Classification: Classification is a fundamental
 task in SAR image interpretation that evaluates the
 VLM's basic visual understanding through target
 category discrimination.

(2) Fine-Grained Description: The fine-grained description task focuses on both target category identification and geometric attribute analysis in SAR imagery. Beyond basic classification, it evaluates the VLM's capability to extract detailed morphological features and spatial orientations, demonstrating the model's proficiency in reasoning about SAR-specific spatial-geometric relationships.

(3) Instance Counting: This task requires accurate counting of multiple SAR targets while extracting their spatial coordinates and orientation information. The key challenge lies in preventing double-counting errors, particularly in complex scenes where multiple targets overlap. The model must maintain robust counting performance while handling various target densities and background complexities.

(4) Spatial Grounding: This task challenges the 240 model to interpret and reason about complex spatial 241 relationships between multiple targets in SAR imagery, including their relative positions, distances, 243 and directional relationships. The key challenge 244 lies in accurately understanding and describing di-245 verse spatial configurations, especially in scenes 246 with multiple interacting objects and varying spa-247 tial layouts. The model must demonstrate precise 248 spatial reasoning abilities while handling complex 249 multi-target scenarios and maintaining consistent performance across different scene compositions. (5) Cross-Modal Identification: Given specified 252

spatial coordinates, the VLM infers target attributes and generates comprehensive descriptions (size, morphology, direction, distance). This task examines the model's ability to fuse and reason about multimodal information in SAR interpretation.

(6) **Referring:** This reverse-reasoning task challenges the model to locate specific instances in SAR images from textual descriptions. The key challenge lies in bridging semantic-visual gaps while accurately determining target spatial orientations, requiring robust cross-modal reasoning capabilities across varied scene configurations.

3.1.3 Task-Oriented Data Generation

Based on the characteristics of the six tasks, this study designs a multimodal dialogue data generation scheme. The specific rules and implementation logic are as follows, with detailed templates provided in the Appendix A.4:

(0) Dataset Definitions

Our dataset adopts a unified representation scheme across all visual-language tasks to ensure consistency and interpretability. The spatial information is uniformly encoded using the bounding box format { $<x_1><y_1><x_2><y_2>$ }, where (x_1,y_1) and (x_2,y_2) denote the top-left and bottom-right coordinates respectively. Spatial relationships are structured through a standard 3×3 grid system (consisting of top-left, top, top-right, left, middle, right, bottom-left, bottom, bottom-right regions).

To explicitly specify different task requirements, we incorporate task-specific prompts: [count] for Instance Counting Task, [grounding] for Spatial Grounding Task, [identify] for Cross-Modal Identification, and [refer] for Referring Task. These prompts help guide the model's attention to the relevant aspects of each task.

These definitions form the foundational framework for our task formulations and evaluation metrics, enabling systematic assessment of visuallanguage models' capabilities.

(1) Classification Task

The Classification Task assesses the model's SAR image recognition capabilities through 20 distinct question-answer template pairs. Random template combinations enhance data diversity, with standardized notation for multi-target scenarios. (2) Fine-Grained Description Task

Fine-Grained description evaluates the model's structured parsing of satellite imagery through comprehensive quality control. Following our dataset definitions, we filter images below 224×224 pixels

363

364

365

366

367

368

370

371

372

373

374

375

376

377

378

379

381

382

383

384

385

354

355

356

and exclude targets with area ratio R < 1% (Equation 1). Targets with aspect ratios exceeding 10:1 or out-of-bounds coordinates are removed. Size descriptions are categorized using area-ratio thresholds (small: <5%, large: >30%). We construct 40 interaction templates to accommodate multi-target scenarios. The calculation of R is formulated as follows:

305

306

310

311

312

313

314

315

316

319

320

321

324

326

330

332

336

337

338

341

342

343

345

347

353

$$R = \frac{w_{box} \times h_{box}}{W_{img} \times H_{img}} \times 100 \tag{1}$$

where w_{box} and h_{box} denote the width and height of the target bounding box, respectively; W_{img} and H_{imq} represent the width and height of the image. (3) Instance Counting Task

As a fundamental component of our visual reasoning system, this task focuses on evaluating the model's quantitative counting capabilities. We designate 15 question templates with [count] identifiers to specify the task requirements, while utilizing our unified bounding box format for structured output representation. The framework supports extended expressions for multi-instance scenarios through coordinate serialization.

(4) Spatial Grounding Task

Spatial Grounding assessment evaluates the model's proficiency in characterizing structural relationships among multiple target objects. Leveraging our established grid system, we quantify spatial relationships through two primary mechanisms: relative distance metrics (with proximal threshold defined as $(W_{imq} + H_{imq})/8)$ and directional relationships (encompassing horizontal, vertical, and diagonal orientations). The framework incorporates 15 spatial-relationship templates, each prefixed with [grounding] identifiers, conforming to our unified spatial representation scheme.

(5) Cross-Modal Identification

Cross-modal parsing evaluation employs a three-tier feature description system. Spatial positioning utilizes a 3×3 grid partitioning scheme for orientation description. Quantitative classification encompasses five-level size descriptions based on area-ratio R thresholds $(\geq 0.4:$ very large; $\geq 0.25:$ large; $\geq 0.1:$ medium; \geq 0.03:small; <0.03:very small) and morphological analysis through bounding-box aspect ratios $(>1.5:wide-body; 0.67 \le ratio \le 1.5:approximately$ square; <0.67:tall-body).

Feature integration combines spatial-sizemorphological elements into comprehensive target profiles. The system implements 20

differential response templates with a dedicated [identify] instruction identifier and structured output templates.

(6) Referring Task

Referring evaluates cross-modal correlation capabilities between natural language and image regions. Queries follow the pattern "Where is the $\{category\}?$ ", prefixed with [refer] identifiers. The task outputs both precise bounding box coordinates and grid-based orientation descriptions, adhering to our unified spatial representation framework through nested parenthetical notation.

3.2 **Quantitative Analysis of Datasets**

The quantitative analysis in this study focuses on two key dimensions: category distribution and object morphological patterns.

Category	Training	Test
Ship	93,373 (46.98%)	10,741 (44.38%)
Aircraft	40,705 (20.48%)	6,779 (28.01%)
Car	9,561 (4.81%)	1,230 (5.08%)
Tank	24,187 (12.17%)	1,773 (7.33%)
Bridge	27,615 (13.89%)	3,281 (13.56%)
Harbor	3,306 (1.66%)	399 (1.65%)

Table 1: Category Distribution Statistics

(1) Category Distribution Characteristics

As shown in Table 1, the ship category dominates both training and test sets (46.98% and 44.38% respectively), while the harbor category represents less than 2%. A significant distribution shift is observed in the aircraft category, with a 7.53% increase in the test set compared to the training set. Categories such as cars, tanks, and bridges maintain moderate and stable proportions across both sets. This class distribution aligns with realworld SAR imagery characteristics, where certain target types naturally appear more frequently than others due to the inherent nature of SAR remote sensing applications and operational scenarios.

(2) Object Morphology Analysis

This study quantify geometric characteristics using aspect ratio (AR):

$$AR = \frac{h_{box}}{w_{box}} \times 100\% \tag{2}$$

Metric	Training Set	Test Set	Diff-Rate
Mean	1.28	1.26	-0.02
Median	1.062	1.05	-0.017
SD	1.18	0.91	-0.22

Table 2: Aspect Ratio Distribution Comparison

As shown in Table 2, the differences in central tendency between training and test sets are minimal (mean: -0.02, median: -0.017). The test set exhibits a 0.22 lower standard deviation, indicating a more concentrated distribution. The key morphological distribution intervals of targets are illustrated in Appendix A.2.

390

394

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

431

432

The dataset exhibits three distinct morphological categories based on aspect ratio (AR): broadbodied (AR \leq 0.67), nearly square-shaped (0.67 <AR \leq 1.5), and tall-bodied (AR > 1.5). Detailed distribution analysis can be found in Appendix A.2.

4 SARChat-Bench Evaluation Method and Settings

This section details the evaluation methodology of SARChat-Bench, a standardized benchmark suite we designed for comprehensive assessment of VLMs in SAR interpretation. The benchmark covers six fundamental tasks that evaluate the model's core capabilities across information processing, target localization, and semantic understanding, providing multi-dimensional insights into visuallanguage model performance in the SAR domain. The evaluation framework ensures fair and thorough assessment of VLMs' capabilities across different SAR interpretation scenarios.

4.1 Evaluation Metrics

İ

(1) Accuracy: A core metric reflecting model prediction fit, calculated as:

$$Acc = \frac{TP}{TP + FP + FN} \times 100\% \qquad (3)$$

where TP denotes correct positive predictions, FP represents false positive predictions, and FNindicates false negative predictions.

(2) Intersection over Union (IoU): In tasks involving localization, identification, and reference, IoU is a key metric measuring the overlap between predicted and ground-truth bounding boxes (bbox). Higher IoU values indicate greater overlap and better localization performance. All IoU-related calculations in this paper are performed with thresholds of 0.25 and 0.5.

(3) Overall Score Calculation:

430
$$S_m = \sum_{t \in T} a_{m,t} \times \frac{n_t}{\sum_{i \in T} n_i}$$
(4)

Among them, n_t represents the sample size of task t, $a_{m,t}$ denotes the accuracy of model m on

task t, and T is the set of all tasks. The detailed calculation of $a_{m,t}$ for each task can be found in Appendix C.

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

4.2 Assessment Methods

This section elaborates on the specific evaluation method processes for six types of tasks.

(1) **Instance Counting:** Compare predicted and label object counts for single-class evaluation, where counting accuracy is measured by Acc and object localization precision is evaluated using Intersection over Union (IoU).

(2) **Spatial Grounding:** Evaluate spatial accuracy through *IoU*-based bbox matching and abstract position analysis (e.g., "top", "bottom") from natural language descriptions.

(3) Cross-Modal Identification: Calculate *IoU* between predicted and ground-truth bboxes for both single and multiple target scenarios to assess cross-modal matching capability.

(4) **Referring:** Assess referring accuracy through *IoU* metrics in both single-target and multi-target contexts.

(5) Fine-Grained Description: Segment predictions and ground-truth into short phrases, extract category and position information, and compare content sets for detailed description evaluation.

(6) Classification: Compare predicted and ground-truth categories to assess classification accuracy *Acc*.

5 Experiments and Analysis

5.1 Implementation Details

During the fine-tuning stage, the model is trained for 1 epoch with a batch size of 4 (effective batch size 32 with gradient accumulation steps of 4). This study employs LoRA training method with rank=8 and alpha=32 on all linear layers. The learning rate is initialized at 1e-4 with a 0.1 warm-up ratio. All experiments are conducted on 2 NVIDIA A100 GPUs using bfloat16 precision. Training is implemented using the MS-SWIFT framework (Zhao et al., 2024) for efficient distributed training. In the most time-consuming case, the training process for a single model took up to 192 hours to complete.

5.2 Benchmark Evaluation

To verify the effectiveness and practicality of the SARChat-2M dataset, we conducted extensive experiments on SAR image interpretation tasks using sixteen mainstream visual-language models. Our

	Δνα	Ανα	Tasks													
Model	Param	1 Avg	Only	Instance	e Count	Abstract	Spatial	Ground	Cross-l	Modal ID	Multi-ta	rget Ref	Single-ta	rget Ref	Descript	Class
		score	count	IoU=.25	IoU=.5	position	Multi	Single	Multi	Single	IoU=.25	IoU=.5	IoU=.25	IoU=.5	Descript	Class
	8B	92.79	74.14	61.37	52.17	81.25	62.25	87.91	98.84	98.98	37.49	23.46	74.86	60.13	63.43	97.25
IntornVI 25	4B	91.57	72.68	57.54	47.35	83.33	60.89	85.90	98.01	98.76	34.05	18.86	69.92	55.29	58.84	97.27
Intern v L2.5	2B	90.55	71.52	54.11	44.22	50.00	60.81	81.92	97.79	98.63	27.05	13.91	68.50	52.16	56.36	96.69
	1B	88.89	69.87	50.18	39.35	0.00	56.30	82.24	96.98	98.60	22.13	9.94	62.33	44.99	53.30	96.65
DoonSookVI	7B	88.99	20.66	8.49	4.19	64.29	65.32	85.78	98.97	99.05	28.75	13.66	64.34	48.84	51.08	93.23
DeepSeek v L	1.3B	84.01	19.61	4.00	1.32	75.00	60.38	82.00	96.40	97.45	16.11	6.23	53.58	34.28	44.44	47.37
Phi-3.5-vision	4.2B	92.06	72.69	57.48	47.60	62.50	58.85	87.29	98.93	98.59	31.65	17.16	70.95	55.70	59.95	96.42
CIM Edas V	2B	90.20	71.59	51.97	40.37	42.86	59.15	86.33	97.54	98.60	24.15	10.66	65.57	46.46	57.86	97.39
GLM-Euge-v	5B	90.48	73.44	56.30	44.56	75.00	61.38	89.96	96.69	95.96	30.68	15.41	69.36	51.81	61.45	98.02
	7B	91.71	71.00	48.07	35.27	100.00	56.37	93.32	99.27	99.51	19.72	7.66	57.27	38.00	54.65	98.80
mPLUG-Owl3	2B	90.32	67.56	41.56	28.83	75.00	45.65	97.58	98.95	99.42	14.91	5.42	50.46	30.16	41.76	98.31
	1B	89.68	67.03	38.64	24.98	75.00	44.07	97.19	98.72	98.87	11.86	4.12	44.34	24.02	40.16	98.06
Owen? VI	7B	90.76	72.79	58.51	50.24	0.00	64.17	83.87	97.54	99.18	39.11	26.29	70.55	57.04	63.11	97.30
Qwell2-VL	2B	90.27	69.63	53.62	45.47	50.00	59.04	78.49	97.55	99.26	32.60	20.12	65.31	51.53	55.20	96.88
LLaVA-1.5	7B	91.21	71.89	56.89	46.80	57.14	62.70	85.79	97.84	98.42	30.81	15.48	71.89	56.70	61.35	96.90
Yi-VL	6B	84.35	32.62	14.35	9.44	75.00	53.68	72.38	93.63	97.95	7.76	2.69	32.95	16.63	38.15	95.32

Table 3: Performance comparison of different vision-language models

preliminary analysis reveals that these models, despite their strong performance on natural images, struggle significantly with SAR image interpretation without domain-specific fine-tuning, demonstrating the critical importance of SAR-domain adaptation. Specifically, we conducted a detailed before-and-after analysis on InternVL2-8B, which achieved the best performance among all tested models, to quantitatively demonstrate the impact of SAR fine-tuning. The results are presented in Appendix E.

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

503

504

507

509

As shown in Table 3, the evaluated models include recent advances such as Qwen2-VL(Wang et al., 2024), InternVL2.5(Chen et al., 2024), DeepSeekVL(Lu et al., 2024), Phi-3.5vision(Abdin et al., 2024), GLM-Edge-V(GLM et al., 2024), mPLUG-Owl3(Ye et al., 2024), Yi-VL(Young et al., 2024) and LLaVA-1.5(Liu et al., 2023). In experiments of different tasks, each model exhibits diverse characteristics:

(1) **Instance Counting** requires VLMs to identify the number of specific objects in the image. Two leading model families achieve state-of-the-art performance: InternVL2.5 and Qwen-VL2, reaching accuracies of 74.14% and 72.79% respectively. However, the accuracy of most other models fall below 60%, highlighting both the challenging nature of the counting task and the dataset's effectiveness in differentiating model capabilities.

(2) Spatial Grounding evaluates models' capabil-510 ity in spatial localization. For abstract position 511 descriptions, mPLUG-Owl3-7B achieves 100% 512 accuracy, significantly outperforming other mod-513 514 els. The mPLUG-Owl3 family maintains superior performance (>90%) in single-target localization, 515 while other models achieve 80%-85%. However, 516 in multi-target scenarios, most models' accuracy 517 drops to approximately 60%. These results sug-518

gest that accurate multi-target spatial information processing remains a crucial area for future model improvements. 519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

(3) Cross-Modal Identification focuses on the model's ability to build precise connections between visual information and other modal information. In this experiment, the process from image recognition to text description is mainly concerned. The experimental data shows that for both single-target and multi-target tasks, the accuracy rates of most models exceed 90%. Among them, the mPLUG-Owl3-7B model performs the best, with the accuracy rates of single-target and multi-target tasks reaching 99.27% and 99.51% respectively, fully demonstrating the powerful capabilities of large language models in cross-modal identification tasks.

(4) **Referring** challenges models to precisely locate objects in SAR images based on textual descriptions. Our experiments reveal significant performance gaps: models achieve less than 75% accuracy on single-target tasks and below 40% on multi-target scenarios. These results highlight the current limitations in cross-modal alignment, particularly in establishing precise text-to-object correspondences within SAR imagery.

(5) Fine-Grained Description requires the model to provide detailed feature and attribute descriptions of the objects in the image. The experiment shows that the model accuracy rates are in the range of 40%-63%. Among them, models with larger parameters such as Qwen2-VL-7B and InternVL2.5-8B perform outstandingly and can give more detailed and accurate descriptions. In contrast, other models with smaller parameter sizes perform poorly, indicating that the accuracy rate of the fine-grained description task is highly sensitive to the model's parameter size.



Figure 3: **Evaluation examples on SARChat-Bench.** VLM predictions are shown in green/red for correct/incorrect descriptions, with the ground truth in green and the predictions in red boxes. And [Human], [Bot], and [Check] icons denote user input, VLMs response, and standard output, respectively.

(6) Classification evaluates models' ability to categorize images based on their content. According to the table data, regardless of parameter size, series such as InternVL2.5, mPLUG-Owl3, Qwen2-VL, and several other models achieve accuracy rates exceeding 96%. The performance of these VLMs demonstrates competitiveness with traditional vision classification models.

558

562

569

570

571

Summary: We benchmark 16 mainstream VLMs on SARChat-Bench. Model size strongly affects fine-grained description performance but shows little impact on classification. While large models excel in cross-modal and class identify tasks and basic spatial grounding, they struggle with referring, counting, detailed descriptions, and multi-target spatial relationships.

5.3 Edge-side models for SAR Applications

This study has multiple edge-side models (≤5B parameters) trained on SARChat-2M and evaluates
their performances. According to Table 3, it demonstrates that these models exhibit task-specific performance variations, achieving remarkable accuracy in cross-modal identification, while showing potential for improvement in referring tasks.
These models support domain-specific fine-tuning

for rapid task adaptation. After optimization, they can operate efficiently on satellite or ground-edge devices, enabling real-time SAR data processing while reducing dependence on cloud infrastructure and minimizing data transmission costs.

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

5.4 Dialogue Visualizaion

Figure 3 presents examples across six tasks from SARChat-Bench. The model completes these tasks with reasonable performance, and its coordinate predictions align with the ground truth annotations. In the spatial grounding task, the model identifies an additional ship that was not included in the original annotations, suggesting its potential in detecting previously unmarked targets in SAR imagery.

6 Conclusion

This research introduces SARChat-2M, a largescale dataset of two million annotated SAR imagetext pairs, addressing the scarcity of languagevision data in the SAR domain. The accompanying SARChat-Bench provides a systematic evaluation framework for assessing VLMs in SAR interpretation tasks, facilitating domain-specific knowledge integration and accelerating the development of SAR-oriented VLMs.

606 Limitation

Despite the comprehensive scale of SARChat-2M based on SARDet-100K dataset, the inherent annotation inconsistencies across different SAR data 609 sources may lead to potential limitations. The varying annotation quality could result in missing tar-611 gets or imprecise target delineation. Notably, there 612 exist cases where VLMs successfully identify valid 613 targets that were not originally annotated in the 614 dataset, highlighting the annotation completeness 615 challenge in the current benchmark construction. 616

617 Ethics Statement

618 In this study, all SAR datasets and methodologies are used strictly for academic research purposes, 619 adhering to their respective licenses and data usage agreements. While our research aims to advance the fundamental understanding of SAR image interpretation, we acknowledge that these technologies 623 could potentially be applied to military or defenserelated purposes. We emphasize that the respon-625 sible application of such technologies is crucial, and their deployment should strictly comply with relevant regulations and ethical guidelines. The research community should maintain ongoing discussions regarding the dual-use nature of SAR technologies to ensure their development serves benefi-631 cial purposes while minimizing potential misuse.

References

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*. 633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

- Jie Chen, Zhixiang Huang, Runfan Xia, Bocai Wu, Lei Sheng, Long Sun, and Baidong Yao. 2022. Largescale multi-class sar image target detection dataset-1.0 [ol]. *Journal of Radars*, (1).
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.
- Qimin Cheng, Haiyan Huang, Yuan Xu, Yuzhuo Zhou, Huanying Li, and Zhongyuan Wang. 2022. Nwpucaptions dataset and mlca-net for remote sensing image captioning. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–19.
- Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. 2018. Functional map of the world. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 6172–6180.
- Wenliang Dai, Junnan Li, D Li, AMH Tiong, J Zhao, W Wang, B Li, P Fung, and S Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. arxiv 2023. *arXiv preprint arXiv:2305.06500*, 2.
- Yimian Dai, Minrui Zou, Yuxuan Li, Xiang Li, Kang Ni, and Jian Yang. 2024. Denodet: Attention as deformable multi-subspace feature denoising for target detection in sar images. *arXiv preprint arXiv:2406.02833*.
- Alexey Dosovitskiy. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, et al. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*.
- Yuan Hu, Jianlong Yuan, Congcong Wen, Xiaonan Lu, and Xiang Li. 2023. Rsgpt: A remote sensing vision language model and benchmark. *arXiv preprint arXiv:2307.15266*.
- Kartik Kuckreja, Muhammad Sohail Danish, Muzammal Naseer, Abhijit Das, Salman Khan, and Fahad Shahbaz Khan. 2024. Geochat: Grounded large vision-language model for remote sensing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27831–27840.

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.

687

691

694

698

706

712

713

715

716

717

718

719

720

721

722

724

725

727

731

732

733

734

735

736

737

739

- Xiang Li, Congcong Wen, Yuan Hu, Zhenghang Yuan, and Xiao Xiang Zhu. 2024. Vision-language models in remote sensing: Current progress and future trends. *IEEE Geoscience and Remote Sensing Magazine*.
- Xin Lin, Bo Zhang, Fan Wu, Chao Wang, Yali Yang, and Huiqin Chen. 2023. Sived: A sar image dataset for vehicle detection based on rotatable bounding box. *Remote Sensing*, 15(11):2825.
- Fan Liu, Delong Chen, Zhangqingyun Guan, Xiaocong Zhou, Jiale Zhu, Qiaolin Ye, Liyong Fu, and Jun Zhou. 2024. Remoteclip: A vision language foundation model for remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning.
- Yang Long, Gui-Song Xia, Shengyang Li, Wen Yang, Michael Ying Yang, Xiao Xiang Zhu, Liangpei Zhang, and Deren Li. 2021. On creating benchmark dataset for aerial image interpretation: Reviews, guidances, and million-aid. *IEEE Journal of selected topics in applied earth observations and remote sensing*, 14:4205–4230.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. 2024.
 Deepseek-vl: Towards real-world vision-language understanding. *Preprint*, arXiv:2403.05525.
- Xiaoqiang Lu, Binqiang Wang, Xiangtao Zheng, and Xuelong Li. 2017. Exploring models and data for remote sensing image caption generation. *IEEE Transactions on Geoscience and Remote Sensing*, 56(4):2183–2195.
- Bo Qu, Xuelong Li, Dacheng Tao, and Xiaoqiang Lu. 2016a. Deep semantic understanding of high resolution remote sensing image. In 2016 International Conference on Computer, Information and Telecommunication Systems (CITS), pages 1–5.
- Bo Qu, Xuelong Li, Dacheng Tao, and Xiaoqiang Lu. 2016b. Deep semantic understanding of high resolution remote sensing image. In 2016 International conference on computer, information and telecommunication systems (Cits), pages 1–5. IEEE.
- Gencer Sumbul, Marcela Charfuelan, Begüm Demir, and Volker Markl. 2019. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. In *IGARSS 2019-2019 IEEE International Geoscience and Remote Sensing Symposium*, pages 5901–5904. IEEE.

Chao Wang, Rui Ruan, Zhicheng Zhao, Chenglong Li, and Jin Tang. 2023. Category-oriented localization distillation for sar object detection and a unified benchmark. *IEEE Transactions on Geoscience and Remote Sensing*. 740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759 760

761

762

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

784

785

786

787

788

789

790

791

- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Yuanyuan Wang, Chao Wang, Hong Zhang, Yingbo Dong, and Sisi Wei. 2019a. A SAR dataset of ship detection for deep learning under complex backgrounds. *Remote. Sens.*, 11(7):765.
- Yuanyuan Wang, Chao Wang, Hong Zhang, Yingbo Dong, and Sisi Wei. 2019b. A sar dataset of ship detection for deep learning under complex backgrounds. *remote sensing*, 11(7):765.
- Shunjun Wei, Xiangfeng Zeng, Qizhe Qu, Mou Wang, Hao Su, and Jun Shi. 2020. Hrsid: A high-resolution sar images dataset for ship detection and instance segmentation. *IEEE Access*, 8:120234–120254.
- Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2024. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models. *arXiv preprint arXiv:2408.04840*.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, et al. 2024. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*.
- Zhiqiang Yuan, Wenkai Zhang, Kun Fu, Xuan Li, Chubo Deng, Hongqi Wang, and Xian Sun. 2022. Exploring a fine-grained multiscale method for crossmodal remote sensing image retrieval. *arXiv preprint arXiv:2204.09868*.
- Yang Zhan, Zhitong Xiong, and Yuan Yuan. 2023. Rsvg: Exploring data and models for visual grounding on remote sensing data. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–13.
- Tianwen Zhang, Xiaoling Zhang, Jianwei Li, Xiaowo Xu, Baoyou Wang, Xu Zhan, Yanqin Xu, Xiao Ke, Tianjiao Zeng, Hao Su, Israr Ahmad, Dece Pan, Chang Liu, Yue Zhou, Jun Shi, and Shunjun Wei. 2021a. Sar ship detection dataset (ssdd): Official release and comprehensive data analysis. *Remote Sensing*, 13(18).
- Tianwen Zhang, Xiaoling Zhang, Jianwei Li, Xiaowo Xu, Baoyou Wang, Xu Zhan, Yanqin Xu, Xiao Ke, Tianjiao Zeng, Hao Su, et al. 2021b. Sar ship detection dataset (ssdd): Official release and comprehensive data analysis. *Remote Sensing*, 13(18):3690.

Wei Zhang, Miaoxin Cai, Tong Zhang, Yin Zhuang, and Xuerui Mao. 2024. Earthgpt: A universal multimodal large language model for multi-sensor image comprehension in remote sensing domain. *IEEE Transactions on Geoscience and Remote Sensing*.

793

794

795

796

797

798

799

800

801

802

803

804 805

806 807

808 809

- Zilun Zhang, Tiancheng Zhao, Yulong Guo, and Jianwei Yin. 2023. Rs5m: A large scale vision-language dataset for remote sensing vision-language foundation model. *arXiv preprint arXiv:2306.11300*.
- Yuze Zhao, Jintao Huang, Jinghan Hu, Xingjun Wang, Yunlin Mao, Daoze Zhang, Zeyinzi Jiang, Zhikai Wu, Baole Ai, Ang Wang, Wenmeng Zhou, and Yingda Chen. 2024. Swift:a scalable lightweight infrastructure for fine-tuning. *Preprint*, arXiv:2408.05517.
- Wang Zhirui, Kang Yuzhuo, Zeng Xuan, WANG Yuelei, ZHANG Ting, and SUN Xian. 2023. Sar-aircraft-1.0: High-resolution sar aircraft detection and recognition dataset. *Journal of Radars*, 12(4):906–922.

811 Appendix

A Data structure analysis

A.1 Word Frequency Analysis of SARChat-2M

As shown in Figure 4, location words (such as center, middle, top) and target object words (such as ship, aircraft, tank) have the highest occurrence frequencies in SAR image descriptions, and the adjective "small" is the most frequently used descriptive word.



Figure 4: Cloud Map of Word-frequency Distribution



Figure 5: The Proportion Distribution of Samples in the Training Set



Figure 6: The Proportion Distribution of Samples in the Testing Set

A.2 Analysis of Datasets Composition

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

Training Set Category Distribution

As shown in Figure 5, in the SARChat-2M training set, the category distribution shows significant differences. Among them, the "Ship" category has the largest proportion, reaching 46.98%, followed by the "Aircraft" category, with a proportion of 20.48%. These two categories account for the majority of the samples in the training set. It can be seen that the sample distribution in the training set is imbalanced, and the "Ship" and "Aircraft" categories dominate. This may enable the model to learn the features of these two categories more comprehensively during the training process. However, since the samples of other categories are relatively few, the model's ability to learn and generalize their features may be affected to a certain extent.

Test Set Category Distribution

As shown in Figure 6, in the SARChat-2M test set, the "Ship" category has the largest proportion among all categories, reaching 44.38%, and the "Aircraft" category ranks second, with a proportion of 28.01%. The distribution trends of these two major categories in the test set are similar to those in the training set. This indicates that the test set has a certain similarity to the training set in terms of the overall category distribution and can be used to test the model's generalization ability on data with a similar distribution. However, the slight differences in the proportions also remind us to comprehensively consider various factors when evaluating the model's performance.

Cross-modal Shape Distribution Analysis

As shown in Figure 7, in the cross-modal identification shape distribution, the "Roughly Square" shape has the largest proportion, with a quantity of 179,652. This shape has an absolute advantage among all shape categories. This means that in the cross-modal identification task, the number of samples of the "Roughly Square" shape is much larger than that of other shapes. The model may be more sensitive to this shape and tend to identify the target as the "Roughly Square" shape during the recognition process. Therefore, when training and optimizing the model, attention should be paid to improving the recognition ability of other shapes to achieve a more balanced recognition effect.

Distribution Analysis of Morphological Categories

As shown in Table 4, nearly square-shaped morphology dominates both datasets, accounting

813

814

815

816

817

818

Category	Dataset	Total Samples	Mean AR	Median AR	Std Dev	AR Distribution (%)				
		F				$AR \le 0.67$	0.67 <ar≤1.5< td=""><td>AR>1.5</td></ar≤1.5<>	AR>1.5		
Ship	Train	93,342	1.34	1.07	1.24	28.37	39.67	31.96		
	Test	10,738	1.308	1.026	1.10	29.34	39.82	30.84		
Aircraft	Train	40,698	1.074	1.047	0.32	5.85	87.67	6.48		
	Test	6,778	1.08	1.041	0.31	4.56	87.36	8.08		
Car	Train	9,561	1.23	1.08	0.56	13.18	60.07	26.75		
	Test	1,230	1.21	1.07	0.53	12.28	62.11	25.61		
Tank	Train	24,15	1.10	1.00	0.84	1.58	94.29	4.13		
	Test	1,771	1.09	1.00	0.29	1.41	94.36	4.23		
Bridge	Train	27,615	1.56	1.18	1.92	18.38	44.83	36.79		
	Test	3,281	1.568	1.2	1.24	18.01	44.59	37.4		
Harbor	Train	3,306	1.20	1.01	0.72	14.19	67.93	17.88		
	Test	399	1.23	1.01	0.81	15.04	68.42	16.54		

Table 4: Analysis of Aspect Ratio of Different Types of Targets



Figure 7: Morphological distribution

for 39.67% in training and 59.37% in test sets, indicating its prevalence in target morphologies. Broad-bodied shapes maintain stable distributions (18.14% training, 17.72% test), while tall-bodied shapes show a moderate decrease from training (31.96%) to test (22.91%) sets. This distribution diversity enhances the model's generalization capability, though the significant increase in nearly square-shaped samples in the test set demands particular attention during model optimization.

Category-Specific Morphological Patterns

As shown in Table 4, each category displays distinctive morphological characteristics. Bridges exhibit the highest average aspect ratio (1.56) with balanced distribution across all morphologies (18% broad, 45% square, 37% tall). Ships demonstrate diverse shapes (28% broad, 40% square, 32% tall), reflecting their real-world variability. Tanks and aircraft show highly concentrated distributions, with nearly square shapes dominating at 94% and 87% respectively, facilitating efficient model learning. Cars and ports maintain moderate distributions with 60-68% nearly square shapes and balanced

remaining proportions.

Analysis of Task based on Datasets A.3 Composition

Our dataset is designed to support unified multimodal tasks through a comprehensive task taxonomy, as shown in Table 5. It encompasses six primary tasks: Classification, Fine-Grained Description, Instance Counting, Spatial Grounding, Cross-Modal Identification, and Referring. Among these, the first three tasks are target-quantity independent, while Spatial Grounding, Cross-Modal Identification, and Referring are further categorized into single-object and multi-object variants. This systematic organization enables diverse training scenarios and enhances model generalization capabilities.

Task Type	Train	Test
Instance Counting	95493 (5.2%)	11794 (5.2%)
Spatial Grounding	94456 (5.1%)	11608 (5.1%)
Cross - Modal Identification	1423548 (77.5%)	175565 (77.4%)
Referring	95486 (5.2%)	11703 (5.2%)
Fine - Grained Description	46141 (2.5%)	6032 (2.7%)
Classification	81788 (4.5%)	10024 (4.4%)

Table 5: Task type distribution in training and test sets



Figure 8: Train Task Distribution

873

895

897

898

899

900

901

902

903

904

905

906

907

908

909



G

Figure 9: Test Task Distribution

As illustrated in Figure 8, the training set com-911 prises 1,836,912 entries. Cross-Modal Identifica-912 tion dominates with 1,423,548 entries (77.50%), 913 enabling robust cross-modal feature learning. In-914 stance Counting and Referring tasks contain 95,493 915 (5.20%) and 95,486 (5.20%) entries respectively, while Spatial Grounding accounts for 94,456 en-917 tries (5.14%). Fine-Grained Description includes 918 46,141 entries (2.51%), with its relatively limited 919 data volume potentially affecting model performance. The Classification task contains 81,788 921 entries (4.45%). 922

924

926

927

929

931

932

935

936

937

941

943

944

The test set (Figure 9) maintains a parallel distribution across its 226,636 entries. Cross-Modal Identification remains dominant with 175.565 entries (77.47%), followed by Instance Counting (11,704, 5.16%), Referring (11,703, 5.16%), Spatial Grounding (11,608, 5.12%), Classification (10,024, 4.42%), and Fine-Grained Description (6,032, 2.66%). This consistent distribution ensures reliable model evaluation.

The dataset exhibits a well-balanced task distribution that supports comprehensive model training. While the consistency between training and test sets ensures reliable evaluation, two aspects warrant attention: the relatively limited data in Fine-Grained Description tasks may constrain model performance and could benefit from expansion, while the dominant proportion of Cross-Modal Identification data necessitates careful consideration of overfitting during training through appropriate regularization methods.

A.4 **Dataset Examples**

Classification Task: User: "Please identify the target categories in the figure (candidate categories: {total category list})" AI: "The {specific category} targets can be identified in the figure'

Spatial Grounding Task:
User: "[grounding] Can you describe the spatial arrangement of various features in this SAR image?" AI: "In this SAR image, I can identify several objects and their locations: There is a {category-1} in the top right portion {bbox}. There are 2 {category-2} in the bottom left: one in the bottom left portion {bbox}, another in the middle left portion {bbox}. The {category-1} is far from the {category-2}."
Fine-Grained Description Task:
<pre>User: "Please describe the content of the satellite image in detail" AI: "In the figure, {size}{category} can be seen located at {orientation},"+"and there is also {size}{category} located at {orientation}"</pre>
Instance Counting Task:
User: "[count] How many {category} are there in the image?"" AI: "There are 3 instances. {bbox}{bbox}
Cross-Modal Identification Task:
<pre>User: "[identify] What object is located in {bbox}?" AI: "This is a medium-sized {category} located in the top right portion of the image. It appears to be wide in shape."</pre>
Referring Task:

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

User: "[refer] Where is the {category}?" AI: "{bbox} (top right)."

Motivation for the Dual-benchmark B Setup

In tasks such as spatial localization, cross-modal identification, and referring expression comprehension, this study constructs a dual-benchmark of single-target and multi-target data. The core research motivations are as follows:

(1) Evaluation of the Model's Target Discrimination Ability

The single-target scenario aims to test the model's basic recognition ability for independent targets. In contrast, the multi-target scenario focuses on examining the model's ability to separate and select targets in complex environments, especially when multiple targets exhibit similar features. This dual-benchmark design can effectively diagnose the performance differences of the model in scenarios with varying degrees of complexity.

(2) Revelation of Target Association Understanding Issues

In multi-target scenarios, the model usually needs to understand the spatial and semantic relationships between targets. By comparing the performance differences of the model in single-target

Performance	8B	4B	2B	1B
Instance Counting Accuracy	74.14	72.68	71.52	69.87
Instance Counting Accuracy ($IoU = 0.25$)	61.37	57.54	54.11	50.18
Instance Counting Accuracy ($IoU = 0.5$)	52.17	47.35	44.22	39.35
Spatial Grounding Accuracy	62.25	60.89	60.81	56.30
Abstract Location in Spatial Grounding Accuracy	81.25	83.33	50.00	0.00
Spatial Grounding Single Accuracy	87.91	85.90	81.92	82.24
Cross-Modal Identification (Multi) Accuracy	98.84	98.01	97.79	96.98
Cross-Modal Identification (Single) Accuracy	98.98	98.76	98.63	98.60
Referring (Multi) Accuracy (IoU = 0.25)	37.49	34.05	27.05	22.13
Referring (Multi) Accuracy ($IoU = 0.5$)	23.46	18.86	13.91	9.94
Referring (Single) Accuracy ($IoU = 0.25$)	74.86	69.92	68.50	62.33
Referring (Single) Accuracy ($IoU = 0.5$)	60.13	55.29	52.16	44.99
Fine-Grained Description Accuracy	63.43	58.84	56.36	53.30
Classification Accuracy	97.25	97.27	96.69	96.65

Table 6: Performance comparison across different model sizes

and multi-target scenarios, it is possible to evaluate
whether the model truly understands the descriptions of the positional relationships between targets.
This helps to identify the limitations of the model
when dealing with relative position descriptions
such as "the vehicle on the left" and "the tank in
the middle".

981

983

985

986

987

988

991

992

994

997

998

1000

1001

1002

1003

1004

1005

1006

1008

(3) Exposure of Attention Mechanism Defects

In multi-target scenarios, the model is highly prone to problems such as attention divergence or overlap. When there are multiple similar targets in an image, the model may have difficulty accurately locating the specific target described by the user. Through the comparison between single-target and multi-target scenarios, the deficiencies of the model in attention allocation can be clearly demonstrated. (4) Simulation of Real-world Application Scenarios

Real-world applications cover both simple single-target scenarios and complex multi-target environments. The establishment of the dualbenchmark is more in line with real-world usage requirements, providing a more comprehensive dimension for model evaluation and helping to improve the applicability and reliability of the model in actual deployments.

C Task-specific Performance Scoring

To evaluate model performance on each task, the task-specific accuacy is caculate by Formula 5. For each task t, the accuracy score $a_{m,t}$ of model m is computed by averaging the accuracy scores across all subtasks:

$$a_{m,t} = \frac{1}{k} \sum_{i=1}^{k} a_{m,t,i}$$
(5)

where $a_{m,t}$ is the average accuracy of model mon task t, k is the number of subtasks, and $a_{m,t,i}$ is the accuracy of model m on the *i*-th subtask of task t. This approach ensures that each subtask contributes equally to the overall task score. 1011

1012

1035

D The Analysis of Model Size

Based on the data analysis in Table 6, it can be con-1013 cluded that for most task-related metrics, there is 1014 a trend of performance improvement as the model 1015 size increases from 1B to 8B. For example, the 1016 instance-counting accuracy rises from 69.87% to 1017 74.14%, the spatial-grounding accuracy increases 1018 from 56.30% to 62.25%, the fine-grained descrip-1019 tion accuracy goes up from 53.30% to 63.43%, and 1020 the classification accuracy climbs from 96.65% to 1021 97.25%. This indicates that an increase in model 1022 size is beneficial to enhancing the performance of 1023 these tasks. However, the accuracy of abstract loca-1024 tions in the spatial-grounding task shows a unique 1025 trend of change. This metric increases from 0.00% 1026 for the 1B model to 83.33% for the 4B model, but 1027 then decreases to 81.25% for the 8B model, not 1028 increasing monotonically with the model size. Evi-1029 dently, the influence of model size on some specific 1030 tasks follows complex patterns. Therefore, when 1031 selecting a model, it is necessary to comprehen-1032 sively consider the task type and model size to 1033 achieve optimal performance. 1034

E Comparison before and after training

To evaluate the impact of the SARChat-2M training 1036 dataset, we conducted a comparative analysis on 1037 InternVL-2.5-8B-the best-performing model in 1038 SARChat-Bench-before and after fine-tuning. Our 1039 evaluation metrics reveal that without fine-tuning 1040 on SARChat-2M, the model fails to comprehend 1041 most SAR-specific targets, exhibiting near-zero per-1042 formance on tasks involving target interpretation 1043 and description. The only exception is the Instance 1044



Table 7: Performance comparison of InternVL2.5 before and after fine-tuning



Figure 10: **Simple Examples of ship detection in SAR images.** The ships appear as distinct bright spots in these SAR images, making them relatively easy to count even for VLMs without SAR-specific training.

Counting Task, where InternVL2.5-8B achieves a baseline accuracy of 63.71%. This relatively high 1046 1047 performance can be attributed to the prevalence of ship-on-sea samples, where SAR imaging exhibits 1048 distinctive characteristics: the smooth sea surface 1049 creates specular reflection, causing most electromagnetic waves to scatter away from the sensor 1051 direction, resulting in weak returns that appear as 1052 dark areas. Meanwhile, ships' metallic structures 1053 and geometric features (such as dihedral and trihe-1054 dral corner reflectors) generate strong backscatter-1055 ing, concentrating radar waves back to the sensor, 1056 thus appearing as bright spots. As shown in Fig-1057 ure 10, these samples present relatively straightforward recognition scenarios, leading to higher count-1059 ing accuracy scores. The comparative results be-1060 tween the base model and its fine-tuned version on 1061 SARChat-Bench are presented in Table 7. Overall, fine-tuning with the SARChat-2M dataset proves essential for enabling VLMs to interpret SAR im-1064 agery effectively. 1065