

SAR-TEXT: From Imperfect Multimodal Earth Observation to Large-Scale SAR–Language Supervision

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Multimodal Earth observation is rarely complete, synchro-*
002 *nized, or uniformly informative in practice. Optical im-*
003 *agery may be unavailable or unreliable under cloud cover,*
004 *nighttime conditions, or adverse weather, whereas synthetic*
005 *aperture radar (SAR) remains observable but is substan-*
006 *tially harder to interpret semantically. This modality gap*
007 *limits the applicability of modern vision–language models*
008 *in realistic remote sensing pipelines. In this paper, we study*
009 *whether large-scale SAR–language supervision can serve*
010 *as a practical bridge for multimodal representation learn-*
011 *ing under imperfect observations. We present SAR-TEXT,*
012 *a 136,584-pair SAR image–text dataset built from hetero-*
013 *geneous SAR sources using a multi-stage caption genera-*
014 *tion pipeline, including annotation-to-caption conversion,*
015 *segmentation-guided captioning, and rule-guided rewriting*
016 *from optical descriptions. We further adopt a progressive*
017 *transfer strategy that adapts vision–language foundation*
018 *models from natural images to optical remote sensing and*
019 *then to SAR. Experiments on cross-modal retrieval, caption*
020 *generation, and downstream SAR visual question answer-*
021 *ing show that large-scale SAR–language supervision sub-*
022 *stantially improves performance over direct-transfer base-*
023 *lines. Human auditing further indicates that the automati-*
024 *cally generated captions are generally usable at scale,*
025 *while failure cases reveal the main bottlenecks under im-*
026 *perfect semantic supervision. Our results suggest that SAR–*
027 *language alignment is a promising mechanism for robust*
028 *multimodal remote sensing when observations are hetero-*
029 *geneous, incomplete, weakly paired, or only partially ob-*
030 *servable.*

031 1. Introduction

032 Earth observation increasingly relies on multiple sensing
033 modalities, including optical satellites, synthetic aperture
034 radar (SAR), infrared imagery, and geospatial metadata.
035 In principle, combining these modalities enables richer

scene understanding, improved retrieval, stronger founda- 036
tion models, and more reliable downstream decision- 037
making [4, 11, 18, 22, 25, 29]. In practice, however, mul- 038
timodal monitoring is often imperfect. Optical observa- 039
tions may be corrupted by clouds, shadows, illumination 040
changes, or acquisition gaps; sensor pairs may be only 041
loosely aligned across space and time; and some modalities 042
may be absent at inference time. These realities motivate 043
a central question for realistic multimodal Earth observa- 044
tion: *how can semantic interpretability be preserved when* 045
the observation space is incomplete or heterogeneous? 046

SAR is especially relevant in this setting. Because it is an 047
active sensor and is robust to weather and illumination, SAR 048
often remains available when optical imagery is degraded or 049
missing [14, 31]. At the same time, SAR is notoriously dif- 050
ficult to interpret semantically. Its backscattering-based for- 051
mation mechanism yields patterns that differ sharply from 052
the appearance cues familiar from natural and optical re- 053
mote sensing images. As a result, state-of-the-art vision– 054
language models trained on natural or optical imagery do 055
not directly transfer well to SAR, and the lack of large- 056
scale SAR–text supervision remains a major bottleneck 057
[6, 20, 30, 32]. 058

In this paper, we argue that SAR–language supervision 059
provides a useful bridge for multimodal Earth observa- 060
tion under imperfect conditions. Rather than viewing SAR 061
as an isolated modality, we treat SAR–text alignment as 062
a mechanism for preserving semantic accessibility when 063
other modalities are unreliable, weakly paired, or unavail- 064
able. Building on this perspective, this study is centered 065
on two components: (1) *SAR-TEXT*, a large-scale corpus of 066
136,584 SAR image–text pairs; and (2) a progressive trans- 067
fer strategy that adapts vision–language models from natu- 068
ral images to optical remote sensing and finally to SAR. 069

Our starting point is the observation that manual SAR 070
captioning is expensive, hard to scale, and often requires 071
substantial domain expertise [30]. We therefore use a multi- 072
stage caption generation pipeline, called *SAR-Narrator*, 073
to transform object detection labels, semantic segmenta- 074
tion maps, and aligned optical descriptions into natural- 075

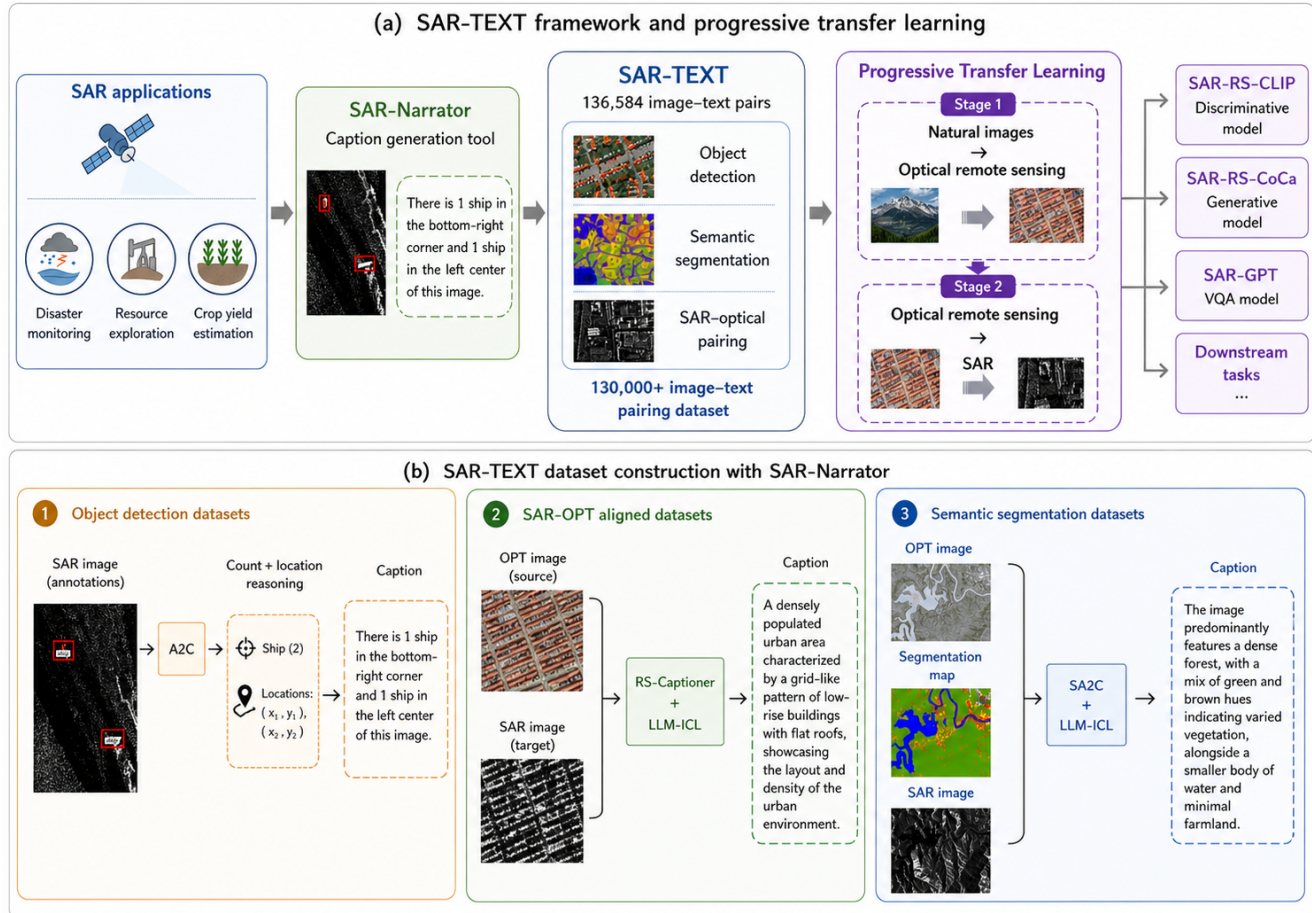


Figure 1. Overview of the proposed SAR–language supervision framework. SAR-Narrator converts heterogeneous SAR supervision sources into SAR-TEXT, which is then used to adapt vision–language models through progressive transfer from natural images to optical remote sensing and finally to SAR.

076 language supervision for SAR images. This yields broad
 077 coverage across object-centric scenes, natural landscapes,
 078 urban areas, and co-registered SAR–optical pairs. We then
 079 study whether such supervision improves retrieval and gener-
 080 ation models in realistic settings where semantic align-
 081 ment must survive strong modality mismatch.

082 Our contributions are fourfold. First, we reframe large-
 083 scale SAR–language supervision as a response to imper-
 084 fect multimodal Earth observation. Second, we present
 085 SAR-Narrator, a practical pipeline for building large-scale
 086 SAR–text corpora from heterogeneous supervision sources.
 087 Third, we show that progressive transfer from natural im-
 088 ages to optical remote sensing and then to SAR improves
 089 both retrieval and caption generation. Fourth, we analyze
 090 human verification and failure modes to characterize what
 091 kinds of semantic errors remain when supervision is gener-
 092 ated automatically at scale, and we further examine SAR-
 093 GPT as a downstream VQA extension of SAR-TEXT.

2. Related Work 094

2.1. Vision–Language Learning for Remote Sensing 095

096 Vision–language learning has become an important direc-
 097 tion in remote sensing, with recent work exploring image-
 098 text alignment, zero-shot recognition, scene understand-
 099 ing, change analysis, and multimodal question answering
 100 [4, 7, 9, 11, 18, 22, 25, 27, 29]. Existing remote sensing
 101 VLMs have primarily focused on optical imagery, where
 102 semantic categories are more accessible and large-scale
 103 image–text supervision is easier to collect. Founda-
 104 tion-style efforts have demonstrated promising results for optical
 105 remote sensing retrieval and transfer, but their assumptions
 106 do not directly carry over to SAR.

107 A key difficulty is that the gap between SAR and opti-
 108 cal imagery is not merely a domain shift; it is a stronger
 109 modality shift rooted in distinct imaging physics. Mod-
 110 els pretrained on natural images often rely on appearance
 111 statistics that do not exist in SAR. This makes naive transfer

Table 1. Source datasets used to construct SAR-TEXT.

| Dataset | Count | Type |
|---------------|---------|------------------|
| MSAR-1.0 | 28,449 | Detection |
| SAR-Ship | 43,819 | Detection |
| OSdataset | 11,245 | SAR-Optical pair |
| HRSID-JPG | 3,642 | Detection |
| QXSLAB_SAROPT | 20,000 | SAR-Optical pair |
| SEN12 | 18,094 | SAR-Optical pair |
| optical_sar | 8,575 | SAR-Optical pair |
| whu-sar-opt | 1,600 | Segmentation |
| SSDD | 1,160 | Detection |
| Total | 136,584 | - |

112 unstable or ineffective, motivating intermediate adaptation
113 through optical remote sensing before final SAR specializa-
114 tion [6, 32].

115 2.2. SAR Image-Text Data

116 Compared with optical remote sensing, SAR image-text re-
117 sources remain limited in scale, diversity, and public acces-
118 sibility [20, 30]. Early attempts at SAR captioning demon-
119 strated feasibility but were typically constrained to narrow
120 object categories or private datasets [30]. More recent mul-
121 timodal remote sensing resources may include some SAR
122 samples, but SAR usually represents only a small portion of
123 the full corpus, and the resulting captions often lack SAR-
124 specific expressiveness [17, 28, 29].

125 This lack of supervision is especially problematic under
126 imperfect multimodal settings. When optical data are miss-
127 ing or unreliable, SAR may become the primary source of
128 evidence. Without language supervision tailored to SAR,
129 downstream multimodal models lose semantic interpretabil-
130 ity precisely when robustness is most needed.

131 2.3. Imperfect Multimodal Observation

132 Real-world multimodal Earth observation often requires
133 monitoring under incomplete, asynchronous, and uncertain
134 observations. In remote sensing, such imperfections arise
135 naturally from weather, seasonality, sensor availability, spa-
136 tial misregistration, and temporal mismatch. While many
137 multimodal pipelines assume clean paired data, practical
138 systems must operate even when pairings are weak or one
139 modality is absent. Existing work on SAR-optical fusion,
140 registration, and joint modeling already highlights this chal-
141 lenge [5, 8, 13, 15, 21]. Our work addresses this setting by
142 using SAR-text alignment as a robust semantic channel that
143 can survive strong heterogeneity and partial pairing.

144 3. Method

145 Our approach consists of two parts: building a large-scale
146 SAR-text corpus with heterogeneous supervision sources,

and adapting vision-language models to SAR through pro-
147 gressive transfer. 148

149 3.1. Design Rationale

150 The goal of SAR-TEXT is not only to increase the number
151 of SAR image-text pairs, but also to provide language su-
152 pervision that is robust to heterogeneous and imperfect ob-
153 servation sources. Different SAR datasets expose different
154 forms of supervision: object detection datasets provide re-
155 liable instance-level labels and coarse spatial locations; se-
156 mantic segmentation datasets provide dense land-cover in-
157 formation; and SAR-optical aligned datasets provide richer
158 scene-level context through paired optical observations. In-
159 stead of forcing all datasets into a single annotation for-
160 mat, SAR-Narrator uses source-specific conversion strate-
161 gies and maps them into a shared natural-language space.

162 This design has two advantages. First, it preserves the
163 strongest signal available from each source. Detection an-
164 notations are most reliable for object counts and positions,
165 segmentation annotations are useful for dominant land-
166 cover categories, and optical captions provide high-level
167 scene descriptions when optical imagery is available. Sec-
168 ond, converting all signals into text enables a unified train-
169 ing interface for vision-language models. As a result, the
170 same SAR-TEXT corpus can support both discriminative
171 alignment, such as SAR-text retrieval, and generative inter-
172 pretation, such as SAR captioning.

173 3.2. SAR-Narrator

174 SAR-Narrator converts available structured or paired infor-
175 mation into natural-language supervision for SAR images.
176 It uses three complementary paths depending on the source
177 dataset.

178 **Detection-to-caption (A2C).** For object detection
179 datasets, we convert category labels, instance counts, and
180 coarse spatial layout into text. The captioning logic uses
181 object counts together with a 3×3 spatial partition of the
182 image, producing simple yet effective captions such as
183 “There is 1 ship in the bottom-right corner of this image”
184 or “There are 14 oil tanks in this image.” This branch is
185 useful for ship, aircraft, bridge, and oil-tank scenes. The
186 A2C branch intentionally uses simple templates rather
187 than long free-form descriptions. This choice reduces
188 hallucination because object detection annotations provide
189 reliable category and location information but do not
190 provide global scene context. For each annotated instance,
191 the image is divided into coarse spatial regions, and the
192 object category, count, and approximate position are
193 verbalized. When multiple objects are present, the caption
194 aggregates repeated categories and describes their relative
195 locations. This produces captions that are concise but
196 strongly grounded in the available annotation.

197 **Segmentation-to-caption (SA2C).** For semantic seg-
198 mentation data, we summarize dominant land-cover cate-
199 gories from pixel-level labels and then rewrite the result-
200 ing description into a more natural caption. This branch
201 captures scene-level semantics such as farmland, city, vil-
202 lage, water, forest, and roads. Instead of exposing raw per-
203 centages directly, the final rewriting prioritizes salient scene
204 elements and improves fluency. The SA2C branch is de-
205 signed for scene-level interpretation. Since segmentation
206 maps provide class proportions rather than natural sentence
207 descriptions, we first identify dominant and secondary land-
208 cover categories and then convert them into a short semantic
209 summary. Very small categories are suppressed to avoid un-
210 stable captions caused by annotation noise. The resulting
211 caption emphasizes salient classes and their approximate
212 dominance, which is particularly useful for natural scenes
213 such as farmland, water bodies, forests, and mixed rural re-
214 gions.

215 **Rule-guided rewriting from optical descriptions.** For
216 co-registered SAR–optical pairs, we generate optical cap-
217 tions with a remote sensing captioner and then rewrite them
218 into SAR-suitable descriptions using an LLM under ex-
219 plicit rules. The rewriting removes color-centric language,
220 avoids speculative phrasing, suppresses camera-related de-
221 tails, and preserves only visually grounded content that
222 plausibly transfers to SAR. A small set of manually written
223 in-context examples stabilizes the rewriting process [3, 10].
224 This rewriting branch is the most expressive but also the
225 most vulnerable to semantic drift. To reduce this risk, we
226 constrain the rewriting process with explicit rules. Color-
227 related descriptions are removed because SAR does not en-
228 code visible color; camera-related phrases are suppressed
229 because they are irrelevant to SAR interpretation; and un-
230 certain visual details are avoided unless they are supported
231 by scene structure. The goal is not to translate optical ap-
232 pearance directly into SAR appearance, but to retain high-
233 level semantics that remain valid across the SAR–optical
234 modality gap.

235 This design reflects the imperfect-observation setting
236 directly. Rather than assuming one clean supervision
237 source, SAR-Narrator accepts heterogeneous inputs—
238 boxes, masks, or paired optical observations—and turns
239 them into language signals tailored to SAR.

240 3.3. SAR-TEXT Dataset Construction

241 We build SAR-TEXT from multiple public SAR datasets
242 spanning three source types: detection datasets, seman-
243 tic segmentation datasets, and SAR–optical registration
244 datasets [1, 5, 8, 13, 15, 16, 19, 26]. These sources cover
245 ships, aircraft, bridges, oil tanks, farmland, water, forest,
246 villages, urban structures, ports, and other scene categories.
247 After cleaning corrupted samples and removing duplicates

or near-duplicates using perceptual hashing [24], we obtain
a final corpus of 136,584 SAR image–text pairs with an av-
erage caption length of 16.3 English words.

251 3.4. Quality Control

252 Because SAR-TEXT is constructed automatically, quality
253 control is essential. We apply several filtering steps before
254 training. First, corrupted images and unreadable files are
255 removed. Second, near-duplicate samples are filtered us-
256 ing perceptual hashing to reduce over-representation of re-
257 peated scenes. Third, captions generated from each source
258 are normalized to reduce formatting inconsistency, redun-
259 dant punctuation, and repeated phrases. Finally, captions
260 that are excessively short, excessively long, or structurally
261 malformed are removed.

262 These filtering steps are particularly important for SAR–
263 language learning. Unlike optical imagery, SAR images of-
264 ten contain speckle noise and ambiguous texture patterns,
265 which can amplify the effect of noisy captions during con-
266 trastive learning. By removing duplicated or malformed
267 samples, the final dataset provides more stable supervision
268 for both retrieval and caption generation.

269 Unlike short class labels, many captions encode both
270 scene semantics and coarse spatial structure. This is valu-
271 able for multimodal representation learning because it sup-
272 ports both discriminative retrieval and generative caption-
273 ing.

274 3.5. Human Verification and Error Profile

275 Automatically generated supervision is useful only if its
276 quality is high enough for learning. We therefore audit a
277 random subset of 500 image–caption pairs. Two annotators
278 with remote sensing experience independently score each
279 pair on object correctness, count correctness, spatial cor-
280 rectness, scene correctness, and fluency using a 1–5 Likert
281 scale, and also provide a binary accept/reject label.

282 The overall acceptance rate is 84.4% (422/500). Aver-
283 age scores are 4.24 for count correctness, 4.27 for fluency,
284 3.96 for spatial correctness, 3.88 for object correctness, and
285 3.76 for scene correctness. These results suggest that the
286 captions are generally usable at scale, though scene-level
287 interpretation remains the most challenging aspect.

288 Failure analysis shows that the dominant errors are
289 wrong scene descriptions (35.9%), hallucination (19.2%),
290 and style bias (17.9%), with fewer cases involving missed
291 or wrong objects. This pattern is informative for imperfect
292 multimodal learning: the easier parts of supervision are lo-
293 cal counts and coarse positions, whereas the harder parts
294 involve global scene interpretation and semantic grounding.

295 The error distribution also reveals that different supervi-
296 sion sources fail in different ways. Detection-derived cap-
297 tions are usually reliable for object categories and counts,
298 but they may oversimplify complex scenes. Segmentation-

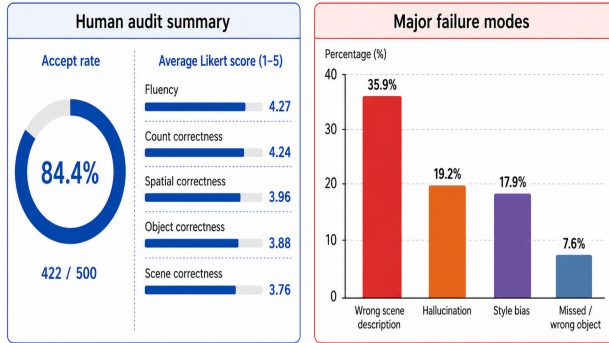


Figure 2. Human verification and failure-mode analysis of SAR-TEXT. The left panel reports the acceptance rate and dimension-wise Likert scores from human auditing, while the right panel summarizes the dominant error categories among rejected samples.

Table 2. Human verification results on 500 sampled SAR-TEXT pairs.

| Metric | Value |
|---------------------|-----------------|
| Sample size | 500 |
| Overall accept rate | 84.4% (422/500) |
| Object correctness | 3.88 ± 0.66 |
| Count correctness | 4.24 ± 0.59 |
| Spatial correctness | 3.96 ± 0.55 |
| Scene correctness | 3.76 ± 0.74 |
| Fluency | 4.27 ± 0.53 |

299 derived captions capture land-cover composition, but they
 300 may miss fine-grained structures and object-level semantics.
 301 SAR-optical rewriting provides richer descriptions,
 302 but it can introduce style bias or retain optical-centric as-
 303 sumptions if the rewriting constraints are insufficient. This
 304 suggests that caption noise in SAR-TEXT is structured
 305 rather than random, which partly explains why the dataset
 306 remains useful for downstream learning despite imperfect
 307 automatic generation.

308 3.6. Progressive Transfer for SAR–Language Align- 309 ment

310 A central difficulty in SAR–language learning is the large
 311 jump from natural-image pretraining to SAR imagery. Di-
 312 rect adaptation may abruptly change the visual embedding
 313 distribution and weaken cross-modal alignment with lan-
 314 guage. We therefore use progressive transfer with an inter-
 315 mediate optical remote sensing stage, following the broader
 316 intuition of transfer learning as staged distribution adapta-
 317 tion [6, 32]:

- 318 1. **Natural/optical adaptation:** start from standard vision–
 319 language backbones and continue pretraining on an op-
 320 tical remote sensing image–text corpus.

2. **SAR adaptation:** fine-tune the adapted models on SAR-
 321 TEXT. 322

323 We apply this strategy to both a CLIP-style retrieval model
 324 and a CoCa-style generative model. The intuition is that
 325 optical remote sensing shares semantic viewpoint and scene
 326 structure with SAR better than natural imagery does, while
 327 still remaining closer to the visual statistics learned during
 328 natural-image pretraining.

4. Experiments 329

4.1. Setup 330

331 We evaluate retrieval on HRSID-test [19] and OSdataset-
 332 512-test, representing object-centric and scene-diverse SAR
 333 scenarios, respectively. For retrieval, we report image-to-
 334 text Recall@K, text-to-image Recall@K, and mean recall.
 335 For caption generation, we report SPICE, BLEU-1/2/3/4,
 336 METEOR, ROUGE-L, and CIDEr.

337 For the retrieval backbone, we use CLIP ViT-L/14 [12].
 338 For the generative backbone, we use CoCa ViT-L/14 [23].
 339 Training follows a full fine-tuning setup on a single RTX
 340 4090 GPU using the OpenCLIP codebase [2].

4.2. Implementation Details 341

342 For retrieval experiments, we fine-tune a CLIP ViT-L/14
 343 backbone using image–text contrastive learning. Each train-
 344 ing batch contains paired SAR images and captions from
 345 SAR-TEXT, and the model is optimized to align corre-
 346 sponding visual and textual embeddings while separating
 347 non-matching pairs. For caption generation, we fine-tune
 348 a CoCa ViT-L/14 backbone to jointly preserve image–text
 349 alignment and autoregressive captioning ability. The same
 350 SAR-TEXT training split is used for both discriminative
 351 and generative models.

352 In the progressive transfer setting, the models are first
 353 adapted to optical remote sensing image–text data and then
 354 further fine-tuned on SAR-TEXT. This two-stage protocol
 355 is kept consistent across CLIP and CoCa to isolate the ef-
 356 fect of progressive transfer. During evaluation, no optical
 357 images are provided to the SAR models; all retrieval and
 358 caption generation results are obtained from SAR images
 359 alone.

360 For the downstream VQA extension, we construct SAR-
 361 VQA instructions from SAR-TEXT and fine-tune a com-
 362 pact multimodal model to obtain SAR-GPT. This evaluation
 363 is not intended to replace the retrieval and captioning bench-
 364 marks, but to test whether SAR-TEXT provides supervision
 365 that can support interactive semantic reasoning over SAR
 366 images.

4.3. Cross-Modal Retrieval Results 367

368 Across both test sets, direct transfer from generic CLIP
 369 performs poorly on SAR. Fine-tuning on SAR-TEXT im-

Table 3. Retrieval results on HRSID-test.

| Model | i2t-R@1 | i2t-R@5 | i2t-R@10 | t2i-R@1 | t2i-R@5 | t2i-R@10 | Mean Recall |
|-------------|-------------|--------------|--------------|-------------|--------------|--------------|--------------|
| CLIP | 0.15 | 0.56 | 1.02 | 0.05 | 0.51 | 0.92 | 0.54 |
| HQRS-CLIP | 0.00 | 0.36 | 0.61 | 0.05 | 0.46 | 1.07 | 0.42 |
| SAR-CLIP | 2.09 | 9.64 | 18.71 | 2.29 | 10.91 | 19.58 | 10.54 |
| SAR-RS-CLIP | 2.65 | 10.91 | 20.50 | 2.55 | 11.78 | 20.75 | 11.52 |

Table 4. Retrieval results on OSdataset-512-test.

| Model | i2t-R@1 | i2t-R@5 | i2t-R@10 | t2i-R@1 | t2i-R@5 | t2i-R@10 | Mean Recall |
|-------------|-------------|--------------|----------|-------------|---------|----------|--------------|
| CLIP | 0.71 | 3.07 | 4.48 | 0.71 | 3.77 | 8.02 | 3.46 |
| HQRS-CLIP | 1.89 | 5.90 | 8.96 | 3.54 | 8.02 | 12.26 | 6.76 |
| SAR-CLIP | 4.48 | 14.86 | 23.82 | 4.72 | 20.99 | 29.72 | 16.43 |
| SAR-RS-CLIP | 5.66 | 16.04 | 23.11 | 5.19 | 20.75 | 28.77 | 16.59 |

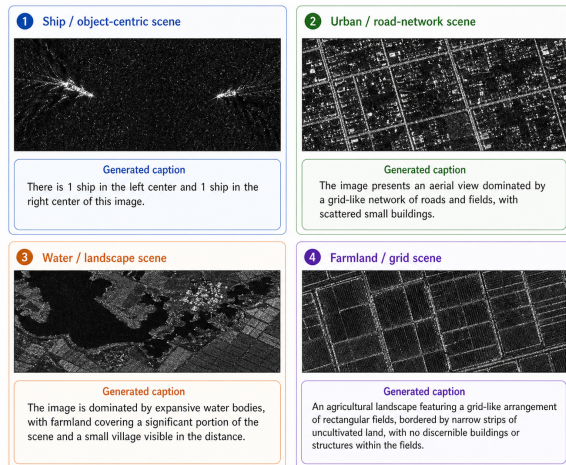


Figure 3. Qualitative captioning examples generated by SAR-RS-CoCa on representative SAR scenes, including object-centric, urban, water-dominated, and farmland-grid scenarios.

proves retrieval substantially, and the progressive strategy improves it further. On HRSID-test, SAR-RS-CLIP reaches a mean recall of 11.52, compared with 0.54 for CLIP and 10.54 for single-stage SAR fine-tuning. On OSdataset-512-test, SAR-RS-CLIP reaches 16.59 mean recall, compared with 3.46 for CLIP and 16.43 for single-stage SAR fine-tuning.

These results support two claims. First, large-scale SAR-language supervision provides effective cross-modal alignment for SAR images. Second, progressive transfer offers a modest but consistent advantage over one-step adaptation, especially on more diverse scenes.

4.4. Caption Generation Results

Generic CoCa performs poorly on SAR, indicating that caption generation does not emerge from natural-image pre-

training alone. Fine-tuning on SAR-TEXT yields large gains, and progressive transfer improves results further. On HRSID-test, SAR-RS-CoCa increases CIDEr from 0.005 for CoCa to 3.186 while also reaching the best BLEU-4, ROUGE-L, and METEOR scores. On OSdataset-512-test, SAR-RS-CoCa similarly yields the strongest overall performance.

These findings show that SAR-language supervision is useful not only for discriminative alignment but also for generative semantic interpretation.

Figure 3 provides qualitative examples from different scene types. In object-centric scenes, the model can identify sparse bright targets and generate captions involving both object counts and relative positions. In urban scenes, the model describes grid-like road structures and scattered buildings, indicating that it has learned scene-level layout cues rather than only isolated bright scatterers. In natural scenes, such as water-dominated and farmland regions, the model captures dominant land-cover patterns and expresses them in coherent language.

At the same time, these examples also illustrate the limitations of SAR captioning. The generated descriptions are generally correct at the coarse semantic level, but they remain less detailed than optical captions and may inherit common phrasing patterns from the training corpus. This observation is consistent with the human verification results, where fluency and count correctness are stronger than fine-grained scene correctness.

The progressive strategy consistently outperforms single-stage mixed training. This supports the idea that adapting first to optical remote sensing and then to SAR provides a smoother transition than jumping directly from natural-image statistics to SAR.

Table 5. Caption generation results on HRSID-test.

| Model | SPICE | B1 | B2 | B3 | B4 | METEOR | ROUGE-L | CIDEr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| CoCa | 0.043 | 0.092 | 0.036 | 0.009 | 0.003 | 0.077 | 0.093 | 0.005 |
| RS-CoCa | 0.059 | 0.081 | 0.039 | 0.014 | 0.007 | 0.088 | 0.121 | 0.001 |
| SAR-CoCa | 0.689 | 0.680 | 0.624 | 0.571 | 0.519 | 0.522 | 0.781 | 2.513 |
| SAR-RS-CoCa | 0.688 | 0.694 | 0.637 | 0.583 | 0.530 | 0.523 | 0.792 | 3.186 |

Table 6. Caption generation results on OSdataset-512-test.

| Model | SPICE | B1 | B2 | B3 | B4 | METEOR | ROUGE-L | CIDEr |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| CoCa | 0.022 | 0.118 | 0.026 | 0.004 | 0.000 | 0.040 | 0.109 | 0.019 |
| RS-CoCa | 0.054 | 0.095 | 0.187 | 0.091 | 0.053 | 0.035 | 0.132 | 0.205 |
| SAR-CoCa | 0.282 | 0.409 | 0.301 | 0.243 | 0.198 | 0.211 | 0.391 | 0.631 |
| SAR-RS-CoCa | 0.298 | 0.420 | 0.319 | 0.263 | 0.219 | 0.228 | 0.410 | 0.665 |

Table 7. Comparison of training strategies on HRSID-test.

| Model | i2t1 | i2t5 | i2t10 | t2i1 | t2i5 | t2i10 | MR |
|--------------------|-------------|--------------|--------------|-------------|--------------|--------------|--------------|
| Single-SAR-RS-CLIP | 2.24 | 8.92 | 16.57 | 1.99 | 9.74 | 17.19 | 9.44 |
| SAR-RS-CLIP | 2.65 | 10.91 | 20.50 | 2.55 | 11.78 | 20.75 | 11.52 |

Table 8. Comparison of training strategies on OSdataset-512-test.

| Model | i2t1 | i2t5 | i2t10 | t2i1 | t2i5 | t2i10 | MR |
|--------------------|-------------|--------------|--------------|-------------|--------------|--------------|--------------|
| Single-SAR-RS-CLIP | 2.83 | 11.32 | 19.34 | 4.01 | 15.80 | 25.71 | 13.17 |
| SAR-RS-CLIP | 5.66 | 16.04 | 23.11 | 5.19 | 20.75 | 28.77 | 16.59 |

418

4.5. Generalization Beyond In-Domain Evaluation

419

420

421

422

423

424

425

426

427

428

429

430

431

To test whether SAR-TEXT transfers beyond benchmarks built with similar annotation procedures, we evaluate on SARLANG-1M-Cap [20]. Because this benchmark does not provide a dedicated retrieval split, the original paper samples 100-image and 500-image subsets for retrieval evaluation and uses the original benchmark for captioning. Fine-tuning on SAR-TEXT substantially improves both retrieval and generation relative to unfine-tuned baselines. The trend is consistent across both retrieval subset sizes, and caption generation metrics improve by roughly twofold after SAR-TEXT fine-tuning. This suggests that SAR-TEXT provides useful supervision that generalizes beyond its own construction pipeline.

432

433

434

435

436

437

438

439

440

This evaluation is important because SAR-TEXT contains captions generated by a specific multi-stage pipeline. A model could in principle overfit to the linguistic style of this pipeline rather than learning transferable SAR-language alignment. Evaluation on SARLANG-1M-Cap provides an additional check against this possibility. Although the benchmark differs in scale and caption style, models fine-tuned on SAR-TEXT still improve over unfine-tuned baselines, suggesting that the learned representations

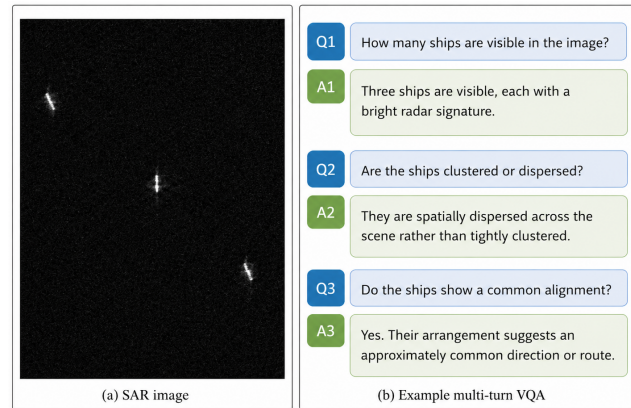


Figure 4. Qualitative example of multi-turn SAR visual question answering with SAR-GPT. Given a SAR image containing several ships, SAR-GPT answers questions about object count, spatial distribution, and approximate alignment, illustrating its ability to support interactive semantic reasoning over SAR imagery.

are not limited to the exact templates used during SAR-TEXT construction. 441

The generalization results also indicate that SAR-language supervision can benefit datasets with different annotation conventions. This is especially relevant for SAR, where existing datasets vary substantially in sensor type, resolution, geographic region, and task definition. A language-aligned representation can provide a common semantic interface across these heterogeneous sources. 442
443
444
445
446
447
448
449

4.6. Downstream Extension: SAR-GPT for SAR VQA 450

Beyond retrieval and caption generation, we also examine whether SAR-TEXT can support instruction-style SAR image understanding. To this end, SAR-VQA is constructed from SAR-TEXT and used to fine-tune a TinyGPT-V-based 451
452
453
454
455

Table 9. Downstream SAR-VQA results of SAR-GPT on representative test sets.

| Split | Model | SPICE | B4 | METEOR | R-L | CIDEr |
|--------|-----------|---------------|---------------|---------------|---------------|---------------|
| HRSID | TinyGPT-V | 0.1789 | 0.0000 | 0.1296 | 0.1975 | 0.2146 |
| | S234-SAR | 0.2959 | 0.1328 | 0.2239 | 0.2975 | 0.9667 |
| | SAR-GPT | 0.3211 | 0.1716 | 0.2350 | 0.3341 | 1.4732 |
| OS-256 | TinyGPT-V | 0.1269 | 0.0215 | 0.1093 | 0.1332 | 0.0483 |
| | S234-SAR | 0.1669 | 0.0730 | 0.1554 | 0.1942 | 0.5513 |
| | SAR-GPT | 0.2974 | 0.2057 | 0.1999 | 0.3276 | 1.7976 |

456 multimodal model, resulting in SAR-GPT. This evaluation
457 is useful because VQA requires a model to map SAR visual
458 evidence into language responses in a more interactive form
459 than standard captioning.

460 As shown in Table 9, SAR-GPT consistently outper-
461 forms TinyGPT-V and S234-SAR on both HRSID and
462 OSdataset-256 SAR-VQA subsets. The gains are partic-
463 ularly clear in CIDEr and BLEU-4, suggesting that SAR-
464 TEXT-derived instruction tuning improves both semantic
465 grounding and response specificity. These results further in-
466 dicate that SAR-TEXT is not limited to image-text match-
467 ing or caption generation, but can also serve as a foundation
468 for downstream SAR-language interaction.

469 5. Discussion Through the Lens of Imperfect 470 Observation

471 The results highlight several implications for robust mul-
472 timodal Earth observation under imperfect sensing condi-
473 tions.

474 **Missing or unreliable optical observations increase**
475 **the value of SAR-text alignment.** When optical imagery
476 is unavailable or less reliable, language-aligned SAR rep-
477 resentations can preserve semantic access to the scene.

478 This is particularly relevant in operational settings where
479 optical imagery is not guaranteed. In such cases, SAR is of-
480 ten the only available source, but its interpretation requires
481 domain expertise. By aligning SAR imagery with language,
482 the model can expose SAR content through a semantic in-
483 terface that is more accessible to downstream retrieval, sum-
484 marization, and human-in-the-loop analysis.

485 **Weak pairing can still be useful if semantic transfor-**
486 **mation is controlled.** The optical-to-SAR rewriting branch
487 shows that imperfect paired supervision can still generate
488 useful language targets when grounded by explicit rules and
489 in-context examples.

490 This observation is important because perfectly synchron-
491 ized, manually described SAR-optical pairs are difficult to
492 obtain at scale. The proposed construction process does not
493 require optical captions to be directly copied into SAR de-
494 scriptions. Instead, it uses optical information as a seman-
495 tic reference and then filters the description through SAR-

aware rewriting rules. This makes weak cross-modal pair- 496
ing a practical source of supervision rather than a strict re- 497
quirement. 498

Noise matters, but structured noise is learnable. Hu- 499
man verification reveals that scene-level semantics remain 500
the main failure point, while counts and spatial phrases are 501
more reliable. This indicates that automatic SAR-language 502
supervision is imperfect but not arbitrary; its errors are 503
structured enough to support effective training. 504

The structured nature of the noise also suggests con- 505
crete directions for improvement. Detection-derived cap- 506
tions could be enriched with contextual scene cues, 507
segmentation-derived captions could incorporate object- 508
level priors, and SAR-optical rewriting could be made more 509
conservative when modality-specific evidence is weak. 510
These improvements would further reduce hallucination 511
while retaining the scalability of automatic annotation. 512

513 6. Limitations and Future Work

This study has several limitations. First, the generated 514
captions are not fully human-authored and may contain 515
style bias, hallucinations, or optical-centric assumptions 516
introduced during rewriting. Second, scene-level seman- 517
tics remain more error-prone than object counts and coarse 518
spatial descriptions, as also reflected by the human audit. 519
Third, although retrieval, caption generation, and SAR- 520
VQA directly evaluate SAR-language alignment, they do 521
not fully cover all imperfect-observation scenarios. Fu- 522
ture work should therefore examine missing-modality rea- 523
soning, asynchronous SAR-optical pairing, temporal mis- 524
match, caption-noise robustness, disaster monitoring under 525
cloud cover, and uncertainty-aware multimodal fusion. Fi- 526
nally, while progressive transfer is empirically effective, a 527
more formal analysis of why staged adaptation mitigates the 528
natural-to-SAR modality gap remains an open direction. 529

530 7. Conclusion

We presented a study of large-scale SAR-language su- 531
pervision for multimodal Earth observation under imper- 532
fect observations. By constructing SAR-TEXT with SAR- 533
Narrator and adapting vision-language models through pro- 534
gressive transfer, we show that SAR-text alignment im- 535
proves retrieval, caption generation, and downstream SAR 536
visual question answering. Human auditing further shows 537
that the generated supervision is generally usable at scale 538
while revealing remaining bottlenecks in scene-level se- 539
mantics and grounding. Overall, our findings suggest that 540
language-aligned SAR representations can serve as a prac- 541
tical semantic bridge when multimodal Earth observation is 542
incomplete, heterogeneous, or weakly paired. 543

544

References

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

- [1] Jie Chen, Zhixiang Huang, Runfan Xia, Bocai Wu, Lei Sheng, Long Sun, and Baidong Yao. Large-scale multi-class sar image target detection dataset-1.0. *Journal of Radars*, 14: 1488, 2022. 4
- [2] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2818–2829, 2023. 5
- [3] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Tianyu Liu, et al. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022. 4
- [4] Sijun Dong, Libo Wang, Bo Du, and Xiaoliang Meng. Changeclip: Remote sensing change detection with multimodal vision-language representation learning. *ISPRS Journal of Photogrammetry and Remote Sensing*, 208:53–69, 2024. 1, 2
- [5] Meiyu Huang, Yao Xu, Lixin Qian, Weili Shi, Yaqin Zhang, Wei Bao, Nan Wang, Xuejiao Liu, and Xueshuang Xiang. The qxs-saropt dataset for deep learning in sar-optical data fusion. *arXiv preprint arXiv:2103.08259*, 2021. 3, 4
- [6] Zhongling Huang, Zongxu Pan, and Bin Lei. What, where, and how to transfer in sar target recognition based on deep cnns. *IEEE Transactions on Geoscience and Remote Sensing*, 58(4):2324–2336, 2019. 1, 3, 5
- [7] Pallavi Jain, Dino Ienco, Roberto Interdonato, Tristan Berchoux, and Diego Marcos. Senclip: Enhancing zero-shot land-use mapping for sentinel-2 with ground-level prompting. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5656–5665. IEEE, 2025. 2
- [8] Xue Li, Guo Zhang, Hao Cui, Shasha Hou, Shun Yao Wang, Xin Li, Yujia Chen, Zhijiang Li, and Li Zhang. Mcanet: A joint semantic segmentation framework of optical and sar images for land use classification. *International Journal of Applied Earth Observation and Geoinformation*, 106: 102638, 2022. 3, 4
- [9] Xiang Li, Congcong Wen, Yuan Hu, and Nan Zhou. Rs-clip: Zero shot remote sensing scene classification via contrastive vision-language supervision. *International Journal of Applied Earth Observation and Geoinformation*, 124:103497, 2023. 2
- [10] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. 4
- [11] Fan Liu, DeLong Chen, Zhangqingyun Guan, Xiaocong Zhou, Jiale Zhu, Qiaolin Ye, Liyong Fu, and Jun Zhou. Remotclip: A vision language foundation model for remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–16, 2024. 1, 2
- [12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. Pmlr, 2021. 5
- [13] Michael Schmitt, Lloyd Haydn Hughes, and Xiao Xiang Zhu. The sen1-2 dataset for deep learning in sar-optical data fusion. *arXiv preprint arXiv:1807.01569*, 2018. 3, 4
- [14] Arsenios Tsokas, Maciej Rysz, Panos M Pardalos, and Kathleen Dipple. Sar data applications in earth observation: An overview. *Expert Systems with Applications*, 205:117342, 2022. 1
- [15] Yuanyuan Wang and Xiao Xiang Zhu. The sarptical dataset for joint analysis of sar and optical image in dense urban area. In *IGARSS 2018-2018 IEEE International Geoscience and Remote Sensing Symposium*, pages 6840–6843. IEEE, 2018. 3, 4
- [16] Yuanyuan Wang, Chao Wang, Hong Zhang, Yingbo Dong, and Sisi Wei. A sar dataset of ship detection for deep learning under complex backgrounds. *remote sensing*, 11(7):765, 2019. 4
- [17] Zhecheng Wang, Rajanie Prabha, Tianyuan Huang, Jiajun Wu, and Ram Rajagopal. Skyscript: A large and semantically diverse vision-language dataset for remote sensing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5805–5813, 2024. 3
- [18] Zhecheng Wang, Rajanie Prabha, Tianyuan Huang, Jiajun Wu, and Ram Rajagopal. Skyscript: A large and semantically diverse vision-language dataset for remote sensing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5805–5813, 2024. 1, 2
- [19] Shunjun Wei, Xiangfeng Zeng, Qizhe Qu, Mou Wang, Hao Su, and Jun Shi. Hrsid: A high-resolution sar images dataset for ship detection and instance segmentation. *Ieee Access*, 8: 120234–120254, 2020. 4, 5
- [20] Yimin Wei, Aoran Xiao, Yexian Ren, Yuting Zhu, Hongruixuan Chen, Junshi Xia, and Naoto Yokoya. Sarlang-1m: A benchmark for vision-language modeling in sar image understanding. *IEEE Transactions on Geoscience and Remote Sensing*, 2026. 1, 3, 7
- [21] Yuming Xiang, Rongshu Tao, Feng Wang, Hongjian You, and Bing Han. Automatic registration of optical and sar images via improved phase congruency model. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 13:5847–5861, 2020. 3
- [22] Kelu Yao, Nuo Xu, Rong Yang, Yingying Xu, Zhuoyan Gao, Titinunt Kitrungrotsakul, Yi Ren, Pu Zhang, Jin Wang, Ning Wei, et al. Falcon: A remote sensing vision-language foundation model. *arXiv preprint arXiv:2503.11070*, 2025. 1, 2
- [23] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022. 5
- [24] Christoph Zauner. Implementation and benchmarking of perceptual image hash functions.(2010). URL http://www.phash.org/docs/pubs/thesis_zauner.pdf, 2010. 4
- [25] Yang Zhan, Zhitong Xiong, and Yuan Yuan. Skyeyegpt: Unifying remote sensing vision-language tasks via instruction

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

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 tuning with large language model. *ISPRS Journal of Pho-*
659 *togrammetry and Remote Sensing*, 221:64–77, 2025. 1, 2
- 660 [26] Tianwen Zhang, Xiaoling Zhang, Jianwei Li, Xiaowo Xu,
661 Baoyou Wang, Xu Zhan, Yanqin Xu, Xiao Ke, Tianjiao
662 Zeng, Hao Su, et al. Sar ship detection dataset (ssdd): Offi-
663 cial release and comprehensive data analysis. *Remote Sens-*
664 *ing*, 13(18):3690, 2021. 4
- 665 [27] Wei Zhang, Miaoxin Cai, Tong Zhang, Yin Zhuang, and
666 Xuerui Mao. Earthgpt: A universal multimodal large lan-
667 guage model for multisensor image comprehension in re-
668 mote sensing domain. *IEEE Transactions on Geoscience and*
669 *Remote Sensing*, 62:1–20, 2024. 2
- 670 [28] Zilun Zhang, Tiancheng Zhao, Yulong Guo, and Jianwei
671 Yin. Rs5m: A large scale vision-language dataset for remote
672 sensing vision-language foundation model. *arXiv preprint*
673 *arXiv:2306.11300*, 2(8), 2023. 3
- 674 [29] Zilun Zhang, Tiancheng Zhao, Yulong Guo, and Jianwei Yin.
675 Rs5m and georsclip: A large scale vision-language dataset
676 and a large vision-language model for remote sensing. *IEEE*
677 *Transactions on Geoscience and Remote Sensing*, 2024. 1,
678 2, 3
- 679 [30] Kai Zhao and Wei Xiong. Exploring data and models in
680 sar ship image captioning. *IEEE Access*, 10:91150–91159,
681 2022. 1, 3
- 682 [31] Xiaobing Zhou, Ni-Bin Chang, and Shusun Li. Applications
683 of sar interferometry in earth and environmental science re-
684 search. *Sensors*, 9(3):1876–1912, 2009. 1
- 685 [32] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi,
686 Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A
687 comprehensive survey on transfer learning. *Proceedings of*
688 *the IEEE*, 109(1):43–76, 2020. 1, 3, 5