

# MVT Appendix Materials

## S1: Land-Use Taxonomy and Hierarchical Mapping

The following specifies the hierarchical mapping between first-level land-use categories and their corresponding second-level subcategories, forming a structured land-cover taxonomy. The taxonomy is constructed according to the national land-use classification standard of the People's Republic of China, as described in Bai-ming Chen *et al.* (2007), and provides a consistent semantic hierarchy for land-use interpretation.

### Label Taxonomy

```
1  LAND_USE_TAXONOMY: Dict[str, Dict] = {
2      "01": {
3          "level1_name": "Cultivated Land",
4          "level2": {
5              "011": {"level2_name": "Paddy Field"},
6              "012": {"level2_name": "Irrigated Land"},
7              "013": {"level2_name": "Dry Land"}
8          }
9      },
10     "02": {
11         "level1_name": "Garden Land",
12         "level2": {
13             "021": {"level2_name": "Orchard"},
14             "022": {"level2_name": "Tea Garden"},
15             "023": {"level2_name": "Other Gardens"}
16         }
17     },
18     "03": {
19         "level1_name": "Forest Land",
20         "level2": {
21             "031": {"level2_name": "Forest"},
22             "032": {"level2_name": "Shrubland"},
23             "033": {"level2_name": "Other Forest Land"}
24         }
25     },
26     "04": {
27         "level1_name": "Grassland",
28         "level2": {
29             "041": {"level2_name": "Natural Grassland"},
30             "042": {"level2_name": "Artificial Grassland"},
31             "043": {"level2_name": "Other Grassland"}
```

```
32     }
33 },
34 "05": {
35     "level1_name": "Commercial Service Land",
36     "level2": {
37         "051": {"level2_name": "Retail Land"},
38         "052": {"level2_name": "Accommodation and Catering Land"},
39         "053": {"level2_name": "Business and Financial Land"},
40         "054": {"level2_name": "Other Commercial Land"}
41     }
42 },
43 "06": {
44     "level1_name": "Industrial, Mining, and Storage Land",
45     "level2": {
46         "061": {"level2_name": "Industrial Land"},
47         "062": {"level2_name": "Mining Land"},
48         "063": {"level2_name": "Storage Land"}
49     }
50 },
51 "07": {
52     "level1_name": "Residential Land",
53     "level2": {
54         "071": {"level2_name": "Urban Residential Land"},
55         "072": {"level2_name": "Rural Homestead Land"}
56     }
57 },
58 "08": {
59     "level1_name": "Public Administration and Public Service Land",
60     "level2": {
61         "081": {"level2_name": "Government Land"},
62         "082": {"level2_name": "Press and Publication Land"},
63         "083": {"level2_name": "Scientific and Educational Land"},
64         "084": {"level2_name": "Medical and Charity Land"},
65         "085": {"level2_name": "Cultural Service Land"},
66         "086": {"level2_name": "Public Facilities Land"},
67         "087": {"level2_name": "Park and Green Space"},
68         "088": {"level2_name": "Scenic and Natural Heritage Land"}
69     }
70 },
71 "09": {
72     "level1_name": "Special Land",
73     "level2": {
74         "091": {"level2_name": "Military Facilities Land"},
75         "092": {"level2_name": "Embassies and Consulates Land"},
76         "093": {"level2_name": "Prison Land"},
77         "094": {"level2_name": "Religious Land"},
78         "095": {"level2_name": "Cemetery Land"}
```

```

79      }
80    },
81  "10": {
82    "level1_name": "Transportation Land",
83    "level2": {
84      "101": {"level2_name": "Railway Land"},
85      "102": {"level2_name": "Road Land"},
86      "103": {"level2_name": "Street Land"},
87      "104": {"level2_name": "Rural Road Land"},
88      "105": {"level2_name": "Airport Land"},
89      "106": {"level2_name": "Port Land"},
90      "107": {"level2_name": "Pipeline Transportation Land"}
91    }
92  },
93  "11": {
94    "level1_name": "Water Bodies and Hydraulic Facility Land",
95    "level2": {
96      "111": {"level2_name": "River Surface"},  

97      "112": {"level2_name": "Lake Surface"},  

98      "113": {"level2_name": "Reservoir Surface"},  

99      "114": {"level2_name": "Pond Surface"},  

100     "115": {"level2_name": "Coastal Tidal Flats"},  

101     "116": {"level2_name": "Inland Tidal Flats"},  

102     "117": {"level2_name": "Ditches"},  

103     "118": {"level2_name": "Hydraulic Construction Land"},  

104     "119": {"level2_name": "Glacier and Permanent Snow"}
105   }
106 },
107 "12": {
108   "level1_name": "Other Land",
109   "level2": {
110     "121": {"level2_name": "Idle Land"},  

111     "122": {"level2_name": "Facility Agricultural Land"},  

112     "123": {"level2_name": "Ridge Land"},  

113     "124": {"level2_name": "Saline-Alkali Land"},  

114     "125": {"level2_name": "Swamp"},  

115     "126": {"level2_name": "Sandy Land"},  

116     "127": {"level2_name": "Bare Land"}
117   }
118 }
119 }
120

```

For example, **Cultivated Land (Level-1 ID: 01)** is a first-level category, which contains three Level-2 subcategories: **Paddy Field (011)**, **Irrigated Land (012)**, and **Dry Land (013)**. Each Level-

2 category inherits the semantic scope of its corresponding Level-1 category, forming a deterministic and interpretable coarse-to-fine land-use hierarchy.

## S2: MLLM Finetuning Prompt

### S2.1 First-Step Preprocessing

#### First-Step Category Mapping

```
1 # Fixed category order -> numeric labels (1..29)
2 CATEGORY_ORDER = [
3     "BareLand", "BaseballField", "Beach", "Bridge", "Center", "Church", "Commercial",
4
5     "DenseResidential", "Desert", "Farmland", "Forest", "Industrial", "Meadow", "MediumRe
sidential",
6
7     "Mountain", "Park", "Parking", "Playground", "Pond", "Port", "RailwayStation", "Resort
", "River",
8     "School", "SparseResidential", "Square", "Stadium", "StorageTanks", "Viaduct"
9 ]
```

Before first-step fine-tuning, we **anonymize all semantic class names** by mapping them to a **fixed, deterministic integer label set [1,29]**. Concretely, we define a constant

`CATEGORY_ORDER` and assign labels sequentially according to this order (Specifically, BareLand is labeled as 1, BaseballField as 2, and so on, up to Viaduct as 29 following the given order.). **This anonymization is completed prior to any prompt construction or model training.**

### S2.2 First-Step Finetuning

We apply the following prompt structure as our first-step finetuning period.

#### First-Step Finetuning Prompt

```
1 PROMPT = (
2     "<image>\n"
3     "You are a senior remote-sensing analyst. Carefully examine the remote-
sensing RGB image at multiple scales and infer its numeric label ID (1-29)
strictly from visual evidence.\n"
4     "Internally consider: global spatial layout and landform; geometry, size,
and alignment of man-made structures; texture repetitiveness and granularity;
surface/material cues; linear networks (roads, tracks, embankments,
shorelines); density and relative scale indicators; color/tonal/spectral
```

```

contrasts; cast shadows and illumination; contextual boundaries and
transitions.\n"
5     "Open-set constraint: NEVER state or imply any category or scene name, and
do not verbalize your rationale.\n"
6     "Output format: return ONLY the integer label in [1, 29] as plain text
with no extra words, symbols, or punctuation."
7 )
8
9 INPUT_TXT = (
10    "Return one integer in [1, 29]. Do not include any words, labels, or
explanations."
11 )

```

## S2.3 Second-Step Preprocessing

### Second-Step Category Mapping

```

1 Category Mapping List
2 <original_class_name> → <anonymized_category_id> (number_of_images)
3
4 airfield → category0001 (800)
5 airplane → category0002 (800)
6 airport → category0003 (700)
7 avenue → category0004 (1031)
8 bare_land → category0005 (1137)
9 baseballdiamond → category0006 (800)
10 basketball_court → category0007 (700)
11 beach → category0008 (2118)
12 bridge → category0009 (1136)
13 buildings → category0010 (1117)
14 chaparral → category0011 (800)
15 church → category0012 (700)
16 circular_farmland → category0013 (700)
17 city_road → category0014 (367)
18 coastline → category0015 (653)
19 commercial_area → category0016 (700)
20 container → category0017 (1166)
21 crossroads → category0018 (803)
22 dam → category0019 (434)
23 denseresidential → category0020 (1600)
24 desert → category0021 (2000)
25 farmland → category0022 (2911)
26 forest → category0023 (4457)
27 freeway → category0024 (800)

```

```

28 gamespace → category0025 (1600)
29 golfcourse → category0026 (800)
30 grave → category0027 (465)
31 ground_track_field → category0028 (700)
32 harbor → category0029 (2117)
33 highway → category0030 (306)
34 hirst → category0031 (1046)
35 industrial_area → category0032 (700)
36 intersection → category0033 (1320)
37 island → category0034 (700)
38 lake → category0035 (700)
39 mangrove → category0036 (963)
40 meadow → category0037 (700)
41 mediumresidential → category0038 (800)
42 mobilehomepark → category0039 (100)
43 mountain → category0040 (1664)
44 mountain_road → category0041 (1066)
45 overpass → category0042 (1918)
46 palace → category0043 (700)
47 parkinglot → category0044 (2799)
48 pipeline → category0045 (353)
49 railway → category0046 (1116)
50 railway_station → category0047 (700)
51 rectangular_farmland → category0048 (700)
52 river → category0049 (2071)
53 roundabout → category0050 (700)
54 runway → category0051 (1271)
55 sea → category0052 (1494)
56 snow_mountain → category0053 (967)
57 sparseresidential → category0054 (1600)
58 stadium → category0055 (700)
59 storagetanks → category0056 (2044)
60 stream → category0057 (1109)
61 tennis_court → category0058 (800)
62 terrace → category0059 (700)
63 thermal_power_station → category0060 (700)
64 tower → category0061 (173)
65 turning_circle → category0062 (309)
66 wetland → category0063 (700)
67
68 Total images in step2 finetuning: 67801

```

The anonymized processing of dataset in step 2 is similar in step 1 where we perform a **global category mapping** that deterministically converts all original semantic classes into **anonymous category IDs** (e.g., `airfield → category0001`, ..., `wetland → category0063`), and

the dataset is rebuilt using only these mapped IDs. After this preprocessing, the training data **no longer contains any original class names**.

## S2.4 Second-Step Finetuning

We apply the following prompt structure as our second-step finetuning period.

### Second-Step Finetuning Prompt

```
1 def build_instruction() -> str:
2     return (
3         "You are a remote-sensing imagery expert. <image>\n\n"
4         "You are provided with an image accompanied by detailed segmentation
5         results generated by the SAM model. "
6         "Each segment includes:\n"
7         "- 'bbox': bounding box coordinates [x,y,w,h], approximately locating
8         the segmented object;\n"
9         "- 'area': the area (in pixels) of the segmented object;\n"
10        "- 'counts': pixel-level RLE (Run-Length Encoding) mask precisely
11        defining object boundaries;\n"
12        "- 'score': confidence score assigned by SAM (optional).\n\n"
13
14         "STRICTLY follow these INTERNAL THINKING STEPS (do NOT output your
15         thinking process, these steps are for internal reasoning ONLY):\n\n"
16
17         "Step 1: INTERNALLY analyze the image comprehensively by observing its
18         overall textures, patterns, and visual structure. "
19
20         "Carefully consider the provided segmentation results (bbox, area,
21         pixel-level masks) to identify pixel-level, object-level and region-level
22         features.\n\n"
23
24         "Step 2: Based on your analysis, INTERNALLY select the MOST
25         APPROPRIATE first-level land-use category from the following standardized
26         list:\n"
27
28         "[Cultivated Land, Garden Land, Forest Land, Grassland, Commercial
29         Service Land, Industrial, Mining, and Storage Land, Residential Land, Public
30         Administration and Public Service Land, Special Land, Transportation Land,
31         Water Bodies and Hydraulic Facility Land, Other Land].\n\n"
32
33         "Step 3: INTERNALLY determine the DETAILED subtype within the selected
34         first-level land-use category, precisely guided by the segmentation
35         information provided by SAM.\n\n"
36
37         "Step 4: INTERNALLY form a clear and precise description of this
38         image, explicitly leveraging the SAM segmentation results (bbox locations,
```

```

mask shapes, areas, and pixel-level information) to consolidate your
understanding and classification decision.\n\n"
22
23      "IMPORTANT OUTPUT REQUIREMENT:\n"
24      "You MUST ONLY OUTPUT the final classification category label
corresponding to the detailed land-use subtype (e.g., category0011). "
25      "DO NOT OUTPUT any intermediate reasoning, first-level categories,
subtype names, or image descriptions. "
26      "ONLY OUTPUT a single category ID."
27  )
28

```

The first-level land-use list in the prompt is **not used as label supervision** and does **not introduce semantic leakage**. Our dataset and training targets have already been anonymized to **ID-only category codes**, and the model is explicitly constrained to output **only a single anonymous category ID** (e.g., `category0011`), never any textual class name.

The first-level list is included solely as a **fixed reasoning scaffold** to enforce a **standardized and consistent decision process**, helping the model converge to a unified, deterministic output format under an ID-only objective. Therefore, the learning signal remains strictly **image/segmentation → anonymous ID**, fully decoupled from textual label semantics.

In Step-II training, each sample consists of an image tile paired with the serialized geometric metadata of **all** SAM2-predicted regions in that tile, including bounding boxes, areas, and run-length encoded (RLE) masks.

While the supervised target is an anonymized category ID, **the model performs taxonomy-grounded interpretation conditioned on the collection of region masks**, using them as explicit spatial evidence rather than treating the image as an unstructured global scene.

## S3: LLM-Judge Evaluation

### S3.1 Prompt

#### S3.1.1 Naturalness

##### Role Definition

You are an impartial writing evaluator. Your job is to judge how natural and human-like each DESCRIPTION is.

##### Scope & Constraints

- Work ONLY with the given text, do not use outside knowledge and no fact-checking.

## **Input**

INPUT (DESCRIPTION):

<<DATASET\_JSON>>

---

## **Output Requirements (STRICT)**

- Return exactly five numbers in this order, separated by single spaces:
    - a. Grammar & Syntax
    - b. Lexical Naturalness & Idiomaticity
    - c. Discourse Coherence & Flow
    - d. Style & Register Appropriateness
    - e. Human-likeness vs Machine “Tells”
  - Evaluate in the natural language of the DESCRIPTION.
  - For code-mixed text, It is not a coherent and logical natural language; if unclear, reflect this in Coherence and Human-likeness.
  - Fragments, lists, bullet points, or telegraphic style can be appropriate: assess “Style & Register” relative to whether that style fits a typical description.
  - Penalize machine “tells” : templated phrasing, repetitive n-grams, self-disclosure (e.g., “As an AI...”), over-hedging, rigid slot-filling, or contradictions.
  - Score each DESCRIPTION independently; do not normalize across items.
- 

## **Scoring Rules**

- Each sub-score must be within [0, 5] using 0.5 increments only (allowed set: 0, 0.5, 1, 1.5, ··:, 5).
  - Round to the nearest 0.5 when needed.
- 

## **Formatting Rules**

- Use a dot as the decimal separator.
  - No commas, brackets, units, markdown, explanations, or any extra text.
  - Output must be a single line with single spaces, no leading/trailing spaces.no markdown, no explanations, no extra text.
-

## Special Case

- If the DESCRIPTION is empty or non-textual gibberish, output:  
0 0 0 0

## S3.1.2 Informativeness

### Role Definition

You are an impartial writing evaluator. Rate how detailed, specific, and informative each DESCRIPTION is.

---

### Scope & Constraints

- Work ONLY with the given text, do not use outside knowledge and no fact-checking.
- 

### Input

INPUT (DESCRIPTION):

<<DATASET\_JSON>>

---

### Output Requirements (STRICT)

- Return exactly five numbers in this order, separated by single spaces:
    - Coverage of Key Facets — Breadth of major elements relevant to the described item/topic (within the text itself).
    - Specificity & Quantification — Presence of precise facts (numbers, named entities, model/versions), measurable details, and exact qualifiers.
    - Concreteness & Observability — Tangible attributes a reader could verify from the text alone (actions, measurements, materials, examples).
    - Context, Constraints & Relations — Purpose, conditions, assumptions, dependencies, comparisons, trade-offs, timelines.
    - Relevance & Non-Redundancy — Focus on the topic without filler/boilerplate; minimal repetition of the same points.
- 

### Scoring Rules

- Each sub-score must be within [0, 5] using 0.5 increments only (allowed set: 0, 0.5, 1, 1.5, ..., 5).
  - Round to the nearest 0.5 when needed.
-

## Formatting Rules

- Use a dot as the decimal separator.
  - No commas, brackets, units, markdown, explanations, or any extra text.
  - Output must be a single line with single spaces, no leading/trailing spaces.no markdown, no explanations, no extra text.
- 

## Special Case

- If the DESCRIPTION is empty or non-textual gibberish, output:

0 0 0 0

## S3.2 LLM-Scoring Rule

We also describe the scoring rules in detail in the main text; here we provide a more fine-grained scoring rubric:

### Scoring Rules

For each region-level description, we use GPT-4o as a rubric-based judge to score two axes: Naturalness and Informativeness. Each axis contains five weighted sub-modules (as listed in Table S3). For every sub-module, the judge outputs a score from 0 to 5 using 0.5 increments only (with 0 the worst and 5 the best)

### Score aggregation

For each axis, we compute a weighted score and scale it to [0, 100]: we first normalize each sub-score by dividing by 5, multiply by its weight, sum over the five sub-modules (weights sum to 1), and finally multiply by 100. We report Naturalness and Informativeness separately; when a single overall scalar is needed, we use the average of the two axis scores.

### Determinism and formatting

To reduce randomness and ensure reproducible outputs, we set the sampling temperature to 0 and enforce a single-line fixed numeric output format via structured prompting.

Dimension	Sub-module	Weight
Naturalness	Grammar & Syntax (G)	0.25
	Discourse Coherence & Flow (D)	0.25
	Lexical Naturalness & Idiomaticity (L)	0.20
	Style & Register Appropriateness (S)	0.15
	Human-likeness vs Machine “Tells” (H)	0.15
Informativeness	Coverage of Key Facets (C)	0.25
	Specificity & Quantification (S)	0.25
	Concreteness & Observability (O)	0.20
	Context, Constraints & Relations (X)	0.20
	Relevance & Non-Redundancy (R)	0.10

### S3.3 LLM-as-Judge Calibration, Reliability, Abnormal-Case Analysis and Rules

We employ GPT-4o as an LLM-as-judge to produce scalable, rubric-based scores under a predefined scoring guideline and a strict output schema. All samples are evaluated using identical criteria and the same automatic parsing rules to ensure comparability across methods.

**Reliability and stability of judge scoring.** We evaluate the stability of the LLM-as-judge by checking score robustness under repeated scoring of the same inputs (e.g., measuring score variance and agreement rate across multiple runs). This helps verify that the judge provides sufficiently consistent signals for comparative evaluation.

**Zero-score policy for irregular outputs.** To prevent unparseable or non-compliant responses from introducing bias, we assign a score of 0 to any irregular output, including but not limited to:

- Empty outputs (empty string, whitespace-only responses, or missing valid content);
- Format violations (missing required fields, not parsable into the predefined schema, or structural corruption);
- Clearly invalid content (placeholders, repetitive nonsensical text, or content unrelated to the input).

**Handling overlong/truncated outputs.** For outputs that are truncated due to length limits, we follow a “score what is observable” principle: the judge scores the response based on the available (untruncated) portion. If critical information is missing because it appears after the truncation point, the score is reduced accordingly; truncation itself does not automatically trigger a zero score.

**Abnormal-case diagnosis.** We further analyze zero-score or abnormal cases and summarize the major failure modes, such as empty outputs, format violations, truncation-induced incompleteness, and off-topic/unsupported content. This diagnosis improves interpretability and helps identify systematic issues in generation and evaluation.

## S4: Manual Evaluation

### S4.1 Manual Evaluation Scoring Rule

We also describe the scoring rules in detail in the main text; here we provide a more fine-grained scoring rubric:

Ground-truth Level-1/Level-2 labels for manual scoring are produced via expert re-annotation rather than directly taken from LoveDA names. Specifically, we randomly select 449 LoveDA tiles and evaluate region-text pairs: each “sample” corresponds to one predicted mask/region and its generated tag/description. Five remote-sensing/GIS researchers independently assign the reference Level-1 and Level-2 category for each region following the Chinese Land Use

Classification Category (Reference 17 and Appendix S1) by inspecting the RGB tile with the mask overlay (and the region’s dominant pixels). For descriptions, annotators judge whether the generated text correctly characterizes the same region’s land-use/land-cover attributes and context, using the image–mask evidence as reference (0 = incorrect/mismatched; 0.5 = partially correct but incomplete or mixed; 1 = fully correct and consistent). Disagreements are resolved by majority vote with final adjudication; predicted tags are scored by exact match (0/1) against the consensus Level-1/Level-2 reference, and description scores are averaged across annotators.

<b>Scoring Item</b>	<b>Score Rule</b>
Level-1 Category	0 / 1
Level-2 Category	0 / 1
Description	0 / 0.5 / 1
<b>Scoring Logic</b>	
Level-1 / Level-2 incorrect	0
Description partly correct	0.5
Description fully correct	1
<b>Total Sample Score:</b> 3 points	

## S4.2 Manual Evaluation Protocol, Annotator Details, Fairness Measures and Rules

We complement automatic evaluation with a human assessment protocol designed to be fair, consistent, and domain-appropriate. We recruit five annotators with substantial experience in remote sensing and related fields to score model outputs

Training and guideline standardization. Before formal evaluation, all annotators undergo training to ensure consistent interpretation of the rubric, including walkthroughs of scoring criteria, calibration with representative examples, and discussion of borderline cases. Annotators follow a unified written guideline during evaluation.

Scoring and aggregation. Each sample is scored independently by the annotators according to the same rubric. The final human score is obtained by aggregating individual ratings (e.g., mean across annotators), and disagreements are handled according to the guideline to ensure consistency.

Fairness control. To maintain fairness, all methods are evaluated under identical conditions and with the same scoring criteria. Annotators are instructed to base their judgments strictly on task-relevant quality and compliance, and to avoid using stylistic preferences or non-essential cues as scoring factors. The training phase and calibration rounds are used to reduce inter-annotator variance and improve evaluation reliability.

All five annotators were trained with unified guidelines and calibration examples to align scoring criteria. We further report the inter-annotator agreement for Level-1/Level-2 labeling and the ordinal description scores, showing substantial consistency across annotators and supporting the reliability of the manual evaluation.

## S5: Detailed Setting Illustration in Stage I: SAM2 Segmentation

In Stage I, we adopt SAM2 as a class-agnostic mask proposer, using the official checkpoint as the off-the-shelf baseline and further fine-tuning it on the full OpenEarthMap to obtain an adapted variant. For evaluation on LoveDA (where urban scenes often contain highly fragmented regions and thin, complex structures), we derive instance proxies by decomposing each semantic class mask into per-class connected components, and match predicted masks to these GT instances via one-to-one IoU assignment (unmatched GT = 0) to compute instance mIoU/mDice.