

REMOTE SENSING-ORIENTED WORLD MODEL

**Yuxi Lu¹, Biao Wu¹, Zhidong Li, Kunqi Li¹, Chenya Huang¹
Huacan Wang², Qizhen Lan³, Ronghao Chen⁴, Ling Chen¹, Bin Liang¹**

¹University of Technology Sydney (UTS), Sydney, Australia

²University of Chinese Academy of Sciences (UCAS), Beijing, China

³University of Alabama at Birmingham, Birmingham, USA

⁴Peking University, Beijing, China

ABSTRACT

World models have shown potential in artificial intelligence by predicting and reasoning about world states beyond direct observations. However, existing approaches are predominantly evaluated in synthetic environments or constrained scene settings, limiting their validation in real-world contexts with broad spatial coverage and complex semantics. Meanwhile, remote sensing applications urgently require spatial reasoning capabilities for disaster response and urban planning. This paper bridges these gaps by introducing the first framework for world modeling in remote sensing. We formulate remote sensing world modeling as direction-conditioned spatial extrapolation, where models generate semantically consistent adjacent image tiles given a central observation and directional instruction. To enable rigorous evaluation, we develop RSWISE (Remote Sensing World-Image Spatial Evaluation), a benchmark containing 1,600 evaluation tasks across four scenarios: general, flood, urban, and rural. RSWISE combines visual fidelity assessment with instruction compliance evaluation using GPT-4o as a semantic judge, ensuring models genuinely perform spatial reasoning rather than simple replication. Afterwards, we present RemoteBAGEL, a unified multimodal model fine-tuned on remote sensing data for spatial extrapolation tasks. Extensive experiments demonstrate that RemoteBAGEL consistently outperforms state-of-the-art baselines on RSWISE.

1 INTRODUCTION

World models have emerged as a frontier in artificial intelligence, showing promise across diverse applications such as robotic navigation (Wu et al., 2023) and autonomous driving (Guan et al., 2025). These models aim to learn the underlying dynamics of environments from limited observations and to predict or reason about unobserved states (Ding et al., 2025). However, most world model studies remain confined to synthetic simulators or constrained scene settings. Synthetic settings lack the complexity and uncertainty of real environments, while constrained scene settings fail to capture reasoning over large spatial structures. As a result, the real-world effectiveness of current world models in spatial reasoning remains largely untested.

Remote sensing provides a uniquely powerful testbed for world models. Satellite and aerial imagery naturally encode "world-level" structures such as urban road networks (Yu & Fang, 2023), river systems (Tomsett & Leyland, 2019), agricultural mosaics (Khanal et al., 2020), and forest landscapes (Fassnacht et al., 2024). At the same time, high-impact applications—including flood prediction for disaster response (Nguyen et al., 2024) and infrastructure forecasting in urban planning (Wellmann et al., 2020)—require reasoning beyond directly observed regions. Yet, much of remote sensing research has focused on recognition tasks such as classification (Li et al., 2022; Temenos et al., 2023) and semantic segmentation (Sun et al., 2020; Zhang et al., 2023a), leaving the potential of world modeling in this domain unexplored.

This paper bridges these gaps by introducing the first framework for world modeling in remote sensing. We formulate remote sensing world modeling as direction-conditioned spatial extrapolation (defined in the image-grid frame with up, down, left, and right, rather than geographic cardinal

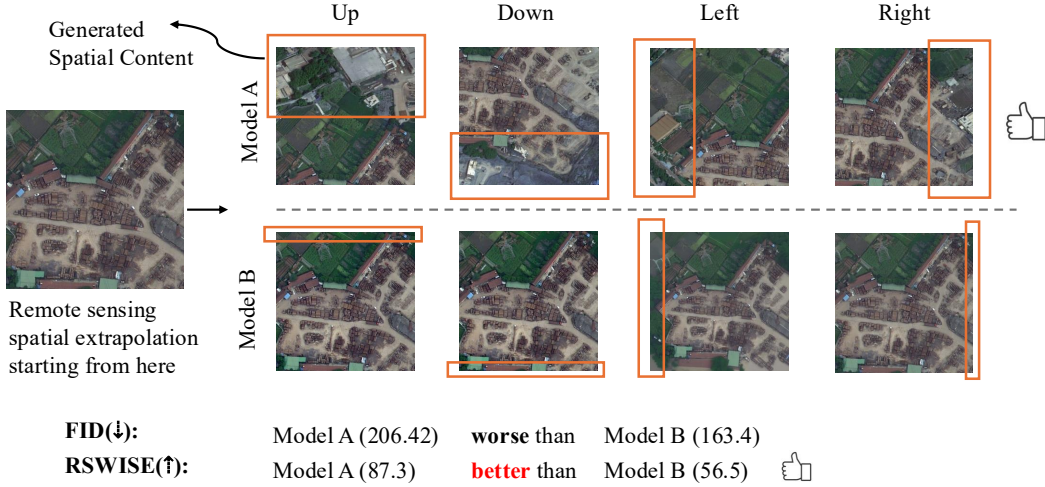


Figure 1: Given a central input tile, models generate direction-conditioned continuations in four directions. Model A produces richer spatial extrapolations (orange boxes), whereas Model B introduces little new content but still attains a lower FID score (163.4 vs. 206.42), which is misaligned with the world modeling objective. The proposed RSWISE metric (87.3 vs. 56.5) accounts for spatial extrapolation consistency, yielding an evaluation more aligned with the requirements of remote sensing world modeling.

directions), where models generate semantically consistent adjacent image tiles given a central observation and directional instruction.

However, introducing world models to the remote sensing domain faces a fundamental evaluation challenge rooted in the limitations of existing assessment paradigms. Current evaluation approaches suffer from critical methodological flaws across two distinct failure modes. Distributional fidelity metrics such as Fréchet Inception Distance (FID) Heusel et al. (2017) measure statistical realism but ignore whether generated tiles follow spatial instructions. As a result, models can obtain deceptively low FID scores by replicating inputs or producing visually plausible but spatially incoherent imagery. Conversely, large language model-based semantic evaluations such as World Knowledge-Informed Semantic Evaluation (WISE) (Niu et al., 2025) capture directional compliance but lack grounding in distributional realism, often rewarding spatially consistent yet environmentally implausible generations.

To overcome these limitations, we introduce RSWISE, the first evaluation framework designed for remote sensing world models. RSWISE integrates distributional fidelity with spatial reasoning consistency through a dual-dimension approach. Specifically, it employs FID to ensure adherence to real-world satellite statistics and leverages GPT-4o to assess whether generated tiles reveal novel yet geographically plausible regions aligned with directional prompts. As shown in Figure 1, RSWISE better reflects geospatial alignment than FID alone, providing a principled basis for fair comparison and progress tracking.

Afterwards, we present RemoteBAGEL, the first remote sensing world model designed for direction-conditioned spatial extrapolation. Instead of producing visually plausible but spatially incoherent completions, RemoteBAGEL explicitly couples generation with spatial reasoning requirements. It is built around two components: (1) a trajectory-based data construction pipeline that transforms raw satellite imagery into instruction-conditioned continuation tasks, and (2) a reconstruction-driven training framework that enforces geographic continuity and semantic coherence during spatial extrapolation. Extensive experiments demonstrate that RemoteBAGEL consistently outperforms state-of-the-art baselines on RSWISE. In summary, our contributions are threefold:

- We propose a novel problem formulation for remote sensing world modeling as direction-conditioned spatial extrapolation tasks;

Benchmark	#Examples	Application Context	Evaluation Metric	Scen. Div.	Geo. Sem.
VBench Ji et al. (2024)	800	Video	FID / Human	✗	✗
WorldModelBench Li et al. (2025)	350	Games / Video	Task-specific	✓	✗
ChronoMagic-Bench Yuan et al. (2024)	1649	Video	Temporal Consistency	✗	✗
TC-Bench Feng et al. (2024)	150	Video	FID	✗	✗
RSWISE (Ours)	1600	Remote Sensing	FID + GPT semantic	✓	✓

Table 1: Comparison of benchmarks. Column headers are abbreviated for readability: Scen. Div. = *Scenario Diversity*, and Geo. Sem. = *Geospatial Semantics*. Existing benchmarks mainly evaluate temporal prediction in robotics, video, or game settings. In contrast, RSWISE (Remote Sensing World-Image Spatial Evaluation) provides 1,600 evaluation tasks, constructed from 100 images \times 4 scenarios \times 4 directions. It focuses on spatial continuation in remote sensing, leveraging real geospatial imagery and explicitly evaluating semantic continuity of structures such as rivers, roads, and urban–rural transitions.

- We introduce RSWISE, the first comprehensive evaluation framework with dual-dimension metrics and a benchmark of 1,600 evaluation tasks across four representative scenarios;
- We develop RemoteBAGEL, the first specialized world model achieving state-of-the-art performance in remote sensing spatial reasoning tasks.

2 RELATED WORK

2.1 WORLD MODELS AND BENCHMARKS

World models generally divided into two perspectives: models that aim to understand the world by abstracting its underlying mechanisms, and models that aim to predict the future by simulating possible evolutions of the environment (Ding et al., 2025). Early efforts such as *World Models* (Ha & Schmidhuber, 2018) focused on abstracting the external world to gain a deep understanding of its underlying mechanisms, while subsequent work including PlaNet (Hafner et al., 2019) and the Dreamer family (Hafner et al., 2020; 2022) introduced recurrent state-space models (RSSMs) that designed to facilitate forward prediction purely within the latent space. More recent advances extend this principle into generative modeling: transformer-based models such as TransDreamer (Chen et al., 2024a) and Genie (Bruce et al., 2024), diffusion and VAE-driven approaches for scene extrapolation and controllable driving (Wang et al., 2025; Cai et al., 2023), and the JEPa family (Assran et al., 2023) that reframe world models as self-supervised abstraction learners. Despite this diversity, evaluation remains a central challenge: existing benchmarks such as VBench (Ji et al., 2024), ChronoMagic-Bench (Yuan et al., 2024), TC-Bench (Feng et al., 2024), and WorldModelBench (Li et al., 2025) focus on controlled or synthetic scene settings, but they do not explicitly involve spatial continuity or geospatial semantics at the remote sensing scale in real-world contexts. This gap prevents current evaluations from testing how well world models reason over real-world structures such as rivers, roads, or urban–rural transitions.

2.2 REMOTE SENSING MODELS

Remote sensing (RS) acquires Earth observations from satellites and aerial platforms, producing imagery that encodes both spectral variation and large-scale spatial structures across diverse environments (Li et al., 2019). Much of RS research has traditionally focused on recognition-oriented tasks, including land-cover classification (Li et al., 2022; Temenos et al., 2023), semantic segmentation (Sun et al., 2020; Zhang et al., 2023a), object detection (Zhang et al., 2023b; Li et al., 2023; Yu & Ji, 2022), and change detection (Bai et al., 2023; Chen et al., 2024b; Li et al., 2024). Building on this recognition paradigm, recent advances in large-scale pretraining have produced RS foundation models such as SkySense (Guo et al., 2024), AnySat (Astruc et al., 2024), SpectralGPT (Hong et al., 2024), and RemoteCLIP (Liu et al., 2024). These models excel at learning transferable representations and support zero-shot inference across downstream tasks (Huo et al., 2025). However, in contrast to world models, current RS methods seldom attempt spatial continuation or reasoning over large geospatial structures.

3 THE RSWISE BENCHMARK

Design overview. The goal of RSWISE is to provide a comprehensive evaluation framework for remote sensing world models that directly addresses the challenge of spatial reasoning in geospatial contexts. It is built around three components: (1) a unified formulation of directional spatial extrapolation, (2) a multi-scenario dataset capturing diverse conditions, and (3) dual-dimension metrics jointly assessing fidelity and reasoning.

3.1 SPATIAL CONTINUATION SPECIFICATION

Problem formulation. We formalize remote sensing world modeling as a directional spatial extrapolation task. Each evaluation instance is represented by a triplet $(T_{\text{input}}, I_{\text{dir}}, T_{\text{target}})$, where T_{input} denotes the observed central tile of a geographic region, I_{dir} is the directional instruction, and T_{target} is the ground-truth adjacent tile. The objective is to model the conditional distribution

$$p_{\theta}(T_{\text{target}} \mid T_{\text{input}}, I_{\text{dir}}), \quad (1)$$

and generate a tile at inference time via

$$T_{\text{generated}} \sim p_{\theta}(\cdot \mid T_{\text{input}}, I_{\text{dir}}). \quad (2)$$

The requirement is that $T_{\text{generated}}$ achieves distributional fidelity with real satellite imagery while preserving semantic coherence with the specified spatial direction. Importantly, directions are defined in the image-grid coordinate frame (up, down, left, right) rather than cardinal North–South–East–West; our analyses therefore study anisotropy in grid-aligned continuations independent of geographic orientation.

Evaluation axes. To capture the complexity of spatial extrapolation, RSWISE defines three complementary axes: (1) *continuity fidelity*, requiring generated tiles to extend geographic structures across boundaries (e.g., roads, rivers, vegetation patches); (2) *semantic transitions*, requiring plausible changes across heterogeneous regions (e.g., urban to rural, land to water); and (3) *directional consistency*, requiring strict adherence to the instructed direction. These axes ensure that evaluation moves beyond visual plausibility and directly probes spatial reasoning.

3.2 DATASET CURATION

To comprehensively evaluate spatial extrapolation under diverse geospatial conditions, we define four representative scenarios that capture complementary challenges in remote sensing world modeling. The first is the general setting of *geographic extrapolation*, where models extend observed regions into adjacent unseen areas across a broad variety of landscapes. This scenario is designed to incorporate diverse scene types within a single setting, enabling a comprehensive assessment of overall model capability rather than focusing on a specific environment. Beyond this general setting, flood scenarios introduce highly dynamic environmental variations, urban regions emphasize continuity across dense built environments, and rural landscapes highlight consistency within natural and agricultural patterns. These scenarios jointly establish a structured basis for assessing spatial extrapolation across stable, dynamic, structured, and natural contexts.

The benchmark dataset consists of 1,600 curated evaluation instances evenly distributed across these four scenarios: general, flood, urban, and rural. The data are sourced from three publicly available datasets: Sky-SA (Zhu et al., 2025) for general scenes, FloodNet (Rahnemoonfar et al., 2021) for flood events, and LoveDA (Wang et al., 2021) for urban and rural landscapes.

Geospatial scenario taxonomy.

- *General*: diverse landscapes including mountains, forests, coastlines, and mixed terrain, serving as a baseline across varied topographies.
- *Flood*: disaster-response contexts with inundated areas and disrupted land cover, testing robustness under dynamic environmental perturbations.
- *Urban*: dense built environments with road networks and building layouts, challenging models to reason over structured spatial patterns.
- *Rural*: agricultural regions, natural vegetation, and sparse settlements, emphasizing the continuity of natural patterns and land-use transitions.

Data construction pipeline. For each scenario, large satellite images are divided into 3×3 overlapping grids of tiles to preserve boundary consistency and spatial autocorrelation. Start tiles are paired with their four cardinal neighbors (up, down, left, right), yielding evaluation triplets. Directional instructions are standardized into fixed prompts (e.g., “Look at what is below this picture”) to ensure fairness across models. Filtering criteria include cloud cover thresholds, resolution consistency, and temporal alignment. A quality assurance process—combining automated checks for artifacts, manual inspection of geographic coherence, and balanced sampling—produces 400 instances per category.

3.3 RSWISE EVALUATION METRICS

The three evaluation axes—continuity fidelity, semantic transitions, and directional consistency—specify the conceptual dimensions along which spatial extrapolation should be assessed. Consequently, these axes are operationalized in RSWISE via two complementary metrics: *distributional fidelity*, which measures the alignment of generated tiles with the statistical properties of real satellite imagery, and *spatial reasoning*, which assesses whether generated tiles follow the instructed direction of extrapolation. These metrics provide a concrete instantiation of the conceptual framework, ensuring that both realism and reasoning are jointly evaluated.

Distributional Fidelity. We employ FID to quantify how well generated tiles align with the statistical properties of real satellite imagery. For remote sensing applications, FID reflects whether generated content follows scenario-specific geographic distributions such as urban density, vegetation cover, or terrain patterns. To facilitate composite scoring, FID values are globally normalized into $s_{\text{fid}} \in [0, 1]$, standardizing units across scenarios and inverting the metric so that higher values correspond to better performance. This transformation ensures comparability across contexts and allows seamless integration with reasoning-based scores.

Spatial Reasoning. We leverage GPT-4o as an external evaluator to assess whether generated tiles reflect meaningful extrapolation in the instructed direction. Valid outputs include feature continuations (e.g., rivers, roads, mountain ridges), coherent land-use transitions (e.g., urban to suburban, forest to agricultural), and natural boundary progressions (e.g., coastlines, watershed divides). The evaluator assigns scores on a $[0, 10]$ scale based on spatial coherence, directional accuracy, and geographic plausibility, which are then normalized to $s_{\text{spatial}} \in [0, 1]$. This metric explicitly penalizes models that replicate the input texture without introducing new, directionally consistent content.

RSWISE. The final score integrates both fidelity and reasoning via a weighted sum:

$$\text{RSWISE}(m, s) = 100 \cdot \left(w_{\text{spatial}} \cdot s_{\text{spatial}}(m, s) + w_{\text{fid}} \cdot s_{\text{fid}}(m, s) \right), \quad (3)$$

where m denotes the model and s the scenario. We assign greater weight to spatial reasoning while retaining distributional fidelity as a grounding constraint. A representative setting is adopted as the default for RSWISE, with detailed validation and sensitivity analyses reported in the experimental section.

4 REMOTE SENSING-ORIENTED WORLD MODEL

We present **RemoteBAGEL**, a remote sensing world model that performs direction-conditioned spatial extrapolation via action-conditioned tile completion. Our approach has two components: (1) a trajectory-based data construction pipeline that converts unlabeled satellite imagery into action-conditioned continuation tasks, and (2) a reconstruction-centric training objective and architecture that enable controllable spatial extrapolation. We first describe the action-conditioned formulation, then detail the training methodology and the architecture overview.

4.1 ARCHITECTURE OVERVIEW

As illustrated in Figure 2 (c), our architecture employs a unified generative framework where the input tile undergoes feature extraction through a visual encoder, while the directional action is transformed into a learned embedding space. These representations are subsequently integrated

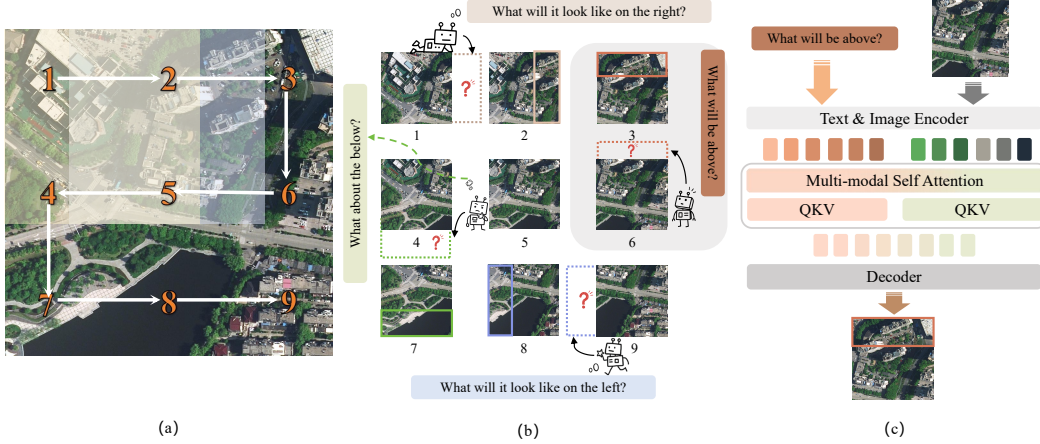


Figure 2: Illustration of the RemoteBAGEL formulation. (a) Large satellite images are partitioned into overlapping 3×3 grids, with example trajectories providing consecutive steps of supervision. (b) Given a central tile and a directional instruction (up, down, left, right), the adjacent tile in the specified direction serves as ground truth, yielding instruction-conditioned triplets. (c) The architecture encodes the input tile and instruction embedding, fuses them via attention, and decodes a continuation tile consistent with the specified direction.

via cross-modal and self-attention mechanisms to capture spatial-semantic dependencies. The fused features are then processed by a generative decoder to synthesize the geographically adjacent tile in the specified direction. This conditioning paradigm enables precise directional control over spatial extrapolation while preserving the structural and semantic coherence inherent in the source imagery.

4.2 ACTION-CONDITIONED DATA CONSTRUCTION

We construct supervision directly from raw satellite imagery without human annotation. As illustrated in Figure 2 (a), large images $X \in \mathbb{R}^{H \times W \times 3}$ are partitioned into overlapping 3×3 grids $\{x_i\}_{i=1}^9$, where overlaps preserve boundary consistency and capture the spatial autocorrelation characteristic of geospatial data. For each central tile x_c , we define a discrete action $a \in \{\text{up, down, left, right}\}$ that specifies a directional move. This yields training triplets

$$(x_c, a, x_{\text{target}}), \quad x_{\text{target}} = \text{adjacent}(x_c, a). \quad (4)$$

in which the adjacent tile in the instructed direction serves as ground truth (Figure 2 (b)). Trajectories (an example route is shown in Figure 2 (a)) provide consecutive steps of supervision, naturally enforcing spatial continuity across tile transitions.

4.3 ACTION-CONDITIONED TRAINING

Given $(x_c, a, x_{\text{target}})$, the model learns a direction-conditioned completion mapping:

$$f_{\theta}(x_c, a) \rightarrow \hat{x}_{\text{target}} \quad (5)$$

and is trained with a reconstruction objective between the prediction and the true neighbor:

$$\mathcal{L}_{\text{recon}} = \mathbb{E}_{(x_c, a)} [\|f_{\theta}(x_c, a) - x_{\text{target}}\|_2^2]. \quad (6)$$

The direction a is encoded as a discrete conditioning token / embedding that modulates generation. Unless otherwise noted, we follow the default loss composition and hyperparameters of the base training recipe, without introducing bespoke auxiliary losses or coefficients. This setup leverages trajectory-based supervision to promote spatial continuity and directional controllability while keeping the objective simple and reproducible.

Model	RSWISE-General	RSWISE-Flood	RSWISE-Urban	RSWISE-Rural	Average
Qwen-Image-Edit	46.9	52.1	56.5	57.2	53.2
FLUX.1-Kontext-Dev	40.0	18.7	43.7	41.8	36.1
Step1X-Edit	51.7	17.3	58.5	55.0	45.6
BAGEL	64.3	64.2	62.3	58.7	62.4
RemoteBAGEL	95.7	78.0	87.3	94.3	88.8

Table 2: Performance of different models on RSWISE. Results are reported across four scenarios (General, Flood, Urban, Rural) and their average. Baselines include Qwen-Image-Edit (Wu et al., 2025), FLUX.1-Kontext-Dev (Labs et al., 2025), Step1X-Edit (Liu et al., 2025), and BAGEL (Deng et al., 2025), while RemoteBAGEL is our proposed method.

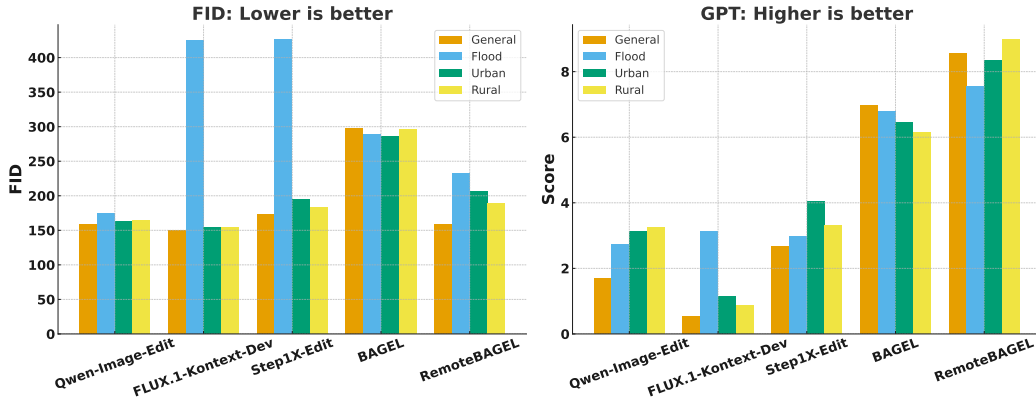


Figure 3: Comparison of models on FID (lower is better) and GPT (higher is better).

5 EXPERIMENT AND RESULTS

5.1 IMPLEMENTATION DETAILS

We fine-tune BAGEL-7B on action-conditioned remote sensing pairs derived from nearly 4,000 images. Each image is partitioned into overlapping 3×3 tiles, resulting in 10,080 direction-labeled training instances aligned with the four RSWISE scenarios. Training is conducted on $4 \times$ H100 (80,GB) GPUs over the course of about 20 hours. For evaluation, we benchmark five models on the RSWISE evaluation set under our protocol, performing inference on $10 \times$ A100 (80,GB) GPUs with roughly 8,000 runs in total, requiring nearly 80 hours.

5.2 MAIN RESULTS

Our evaluation reveals a clear performance hierarchy among the tested models. RemoteBAGEL substantially outperforms all baselines across four benchmark scenarios, achieving near-optimal scores (~ 95) in *general* and *rural* settings. This represents a significant advancement over BAGEL (58-64 points), despite BAGEL’s strong multimodal foundations. The performance gap demonstrates that domain-specific adaptation is crucial for remote sensing tasks, as generic vision-language models struggle to capture the spatial coherence and structural patterns inherent in satellite imagery.

5.3 FID VS. GPT METRICS.

Figure 3 highlights the complementary roles of the two metrics. FID, dominated by texture and artifact statistics, offers a reliable measure of visual fidelity and detail preservation, but is less sensitive to semantic plausibility or directional accuracy. GPT-based evaluation, in contrast, directly probes spatial reasoning by scoring continuity, transitions, and compliance with the instructed direction, yet is less attuned to subtle degradations in low-level image quality. This explains why generic editing baselines such as Qwen-Image-Edit and FLUX.1-Kontext-Dev achieve competitive FID values despite failing to introduce meaningful extrapolated content, as reflected in their low GPT scores.

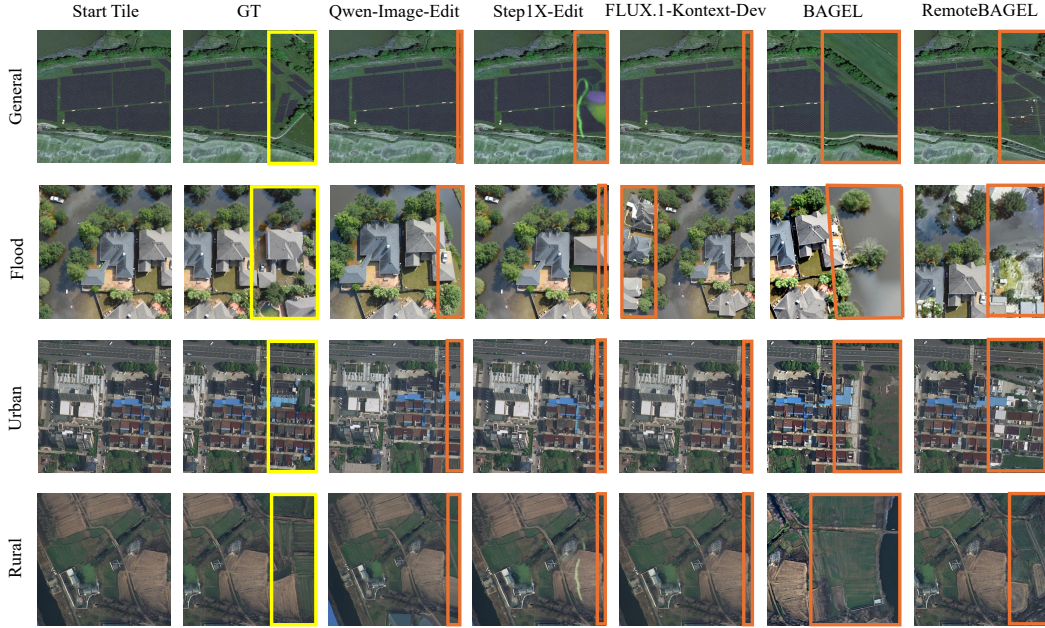


Figure 4: Qualitative comparison of rightward continuations across four benchmark scenarios (general, flood, urban, rural). RemoteBAGEL generates geospatially consistent extrapolations that follow the directional instruction, closely matching the ground truth. By contrast, other models frequently fail to introduce valid extrapolated content or produce semantically inconsistent results (e.g., reversed directions or mismatched structures), leading to large discrepancies from the ground truth.

At the same time, RemoteBAGEL attains the highest GPT scores across all scenarios (up to 8.98 in *rural*) while maintaining competitive FID even in challenging settings such as *flood*. Together, GPT captures the decisive dimension of semantic extrapolation, while FID provides complementary sensitivity to visual quality, sharpening distinctions between strong models like BAGEL and RemoteBAGEL.

5.4 RESULT ANALYSIS

Failure Mode Analysis Qualitative analysis in Figure 4 reveals distinct failure modes across model categories. While BAGEL generates visually plausible outputs, it systematically violates spatial consistency-producing incorrect orientations, missing expected structures (e.g., building clusters), or ignoring geometric regularities in agricultural patterns. Other baselines exhibit more fundamental limitations, generating outputs nearly identical to input tiles with minimal spatial extrapolation, indicating their inability to perform meaningful world modeling. These failures underscore that visual realism alone is insufficient; successful remote sensing continuation requires both semantic understanding and spatial reasoning capabilities that only domain-aware supervision can provide.

Visual Fidelity Analysis Figure 5 (a) compares the FID distributions of BAGEL and RemoteBAGEL across four scenarios. RemoteBAGEL consistently achieves superior visual fidelity with lower FID scores and reduced variance in all settings, though the degree of improvement varies significantly by scenario type. The gains are most pronounced in *general* and *rural* scenarios (FID reductions $> 20\%$), where repetitive agricultural patterns and homogeneous textures benefit substantially from domain-specific training. *Urban* scenarios show moderate but consistent improvements-while geometric regularities in roads and buildings provide structural cues, the higher variance suggests ongoing challenges in capturing fine-grained urban diversity. *Flood* scenarios prove most challenging for both models, exhibiting the smallest improvements due to irregular and dynamic water boundaries that resist systematic pattern learning. These results demonstrate that structured, pattern-rich environments are more suitable to generative modeling than highly variable or transient phenomena.

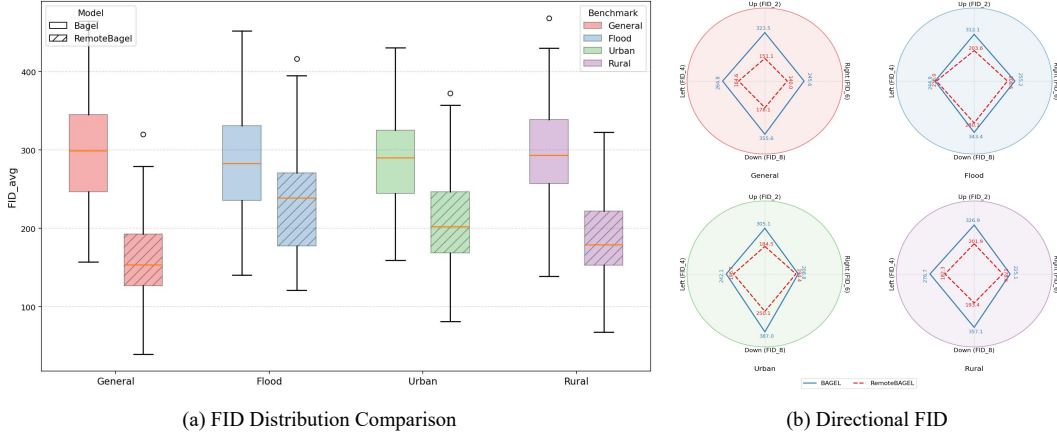


Figure 5: FID-based analysis of BAGEL and RemoteBAGEL. (a) Distributional comparison across four scenarios shows that RemoteBAGEL consistently achieves lower FID, with the largest gains in general and rural scenes and the smallest in flood settings. (b) Directional comparison reveals an anisotropic pattern: left and right continuations are easier to model than upward or downward ones, indicating a directional bias in spatial extrapolation.

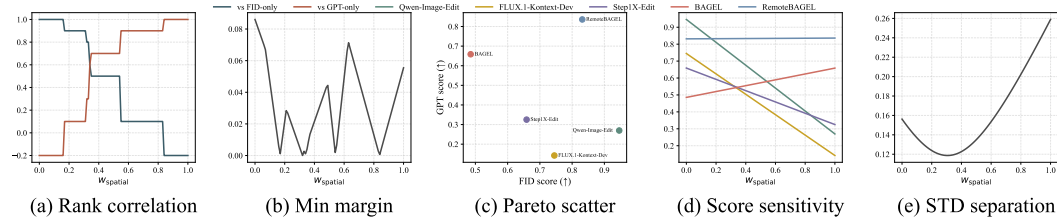


Figure 6: Weight analysis of RSWISE across different criteria. Rankings are stable within $[0.5, 0.7]$, supporting the choice of $w_{\text{spatial}} = 0.6$, $w_{\text{fid}} = 0.4$.

Directional Continuation Analysis Figure 5 (b) reports model performance across the four directional prompts and reveals an anisotropic pattern: continuations along one axis (left-right) tend to achieve lower FID, whereas those along the orthogonal axis (up-down) exhibit higher error rates across all scenarios. This asymmetry does not imply a fixed geographic bias—our tiles are not aligned with cardinal directions—but rather reflects differences in local semantic continuity. Horizontal continuations within the image grid often involve more homogeneous structures such as roads or farmland strips, while vertical continuations more frequently cross heterogeneous transitions such as urban-rural boundaries or land-water interfaces. These results suggest that extrapolation errors are amplified when generation requires bridging semantically diverse regions, highlighting a key challenge for spatial reasoning beyond texture fidelity.

Weight analysis. We validate the RSWISE weighting scheme through a systematic scan of w_{spatial} (Figure 6). The five diagnostics respectively examine rank correlation, min margin, pareto scatter, score sensitivity, and STD separation. Together they show that rankings remain stable and discriminative power is preserved within the interval $[0.5, 0.7]$, supporting the choice of $w_{\text{spatial}} = 0.6$ and $w_{\text{fid}} = 0.4$. Detailed interpretations of each diagnostic are provided in the Appendix.

6 CONCLUSION

Our proposed remote sensing-oriented world model formulates direction-conditioned spatial reasoning and establishes a dedicated benchmark. In the future, this foundation can be expanded by incorporating climate variables, multispectral observations, synthetic aperture radar (SAR), and temporal dynamics, enabling richer capabilities in remote sensing such as cloud removal, weather prediction, and 3D flood visualization.

REFERENCES

- Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture, April 2023.
- Guillaume Astruc, Nicolas Gonthier, Clement Mallet, and Loic Landrieu. AnySat: One Earth Observation Model for Many Resolutions, Scales, and Modalities, 2024.
- Ting Bai, Le Wang, Dameng Yin, Kaimin Sun, Yepi Chen, Wenzhuo Li, and Deren Li. Deep learning for change detection in remote sensing: A review. *Geo-spatial Information Science*, 26(3):262–288, July 2023. ISSN 1009-5020, 1993-5153. doi: 10.1080/10095020.2022.2085633.
- Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge (Jimmy) Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, Yusuf Aytar, Sarah Bechtle, Feryal Behbahani, Stephanie Chan, Nicolas Heess, Lucy Gonzalez, Simon Osindero, Sherjil Ozair, Scott Reed, Jingwei Zhang, Konrad Zolna, Jeff Clune, Nando De Freitas, Satinder Singh, and Tim Rocktäschel. Genie: Generative interactive environments. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*, Vienna, Austria, 2024. JMLR.org.
- Shengqu Cai, Eric Ryan Chan, Songyou Peng, Mohamad Shahbazi, Anton Obukhov, Luc Van Gool, and Gordon Wetzstein. DiffDreamer: Towards Consistent Unsupervised Single-view Scene Extrapolation with Conditional Diffusion Models, March 2023.
- Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. TransDreamer: Reinforcement Learning with Transformer World Models, November 2024a.
- Hongruixuan Chen, Jian Song, Chengxi Han, Junshi Xia, and Naoto Yokoya. ChangeMamba: Remote Sensing Change Detection With Spatiotemporal State Space Model. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–20, 2024b. ISSN 0196-2892, 1558-0644. doi: 10.1109/TGRS.2024.3417253.
- Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified multimodal pretraining, 2025. URL <https://arxiv.org/abs/2505.14683>.
- Jingtao Ding, Yunke Zhang, Yu Shang, Yuheng Zhang, Zefang Zong, Jie Feng, Yuan Yuan, Hongyuan Su, Nian Li, Nicholas Sukiennik, Fengli Xu, and Yong Li. Understanding World or Predicting Future? A Comprehensive Survey of World Models. *ACM Computing Surveys*, pp. 3746449, June 2025. ISSN 0360-0300, 1557-7341. doi: 10.1145/3746449.
- Fabian Ewald Fassnacht, Joanne C White, Michael A Wulder, and Erik Næsset. Remote sensing in forestry: Current challenges, considerations and directions. *Forestry: An International Journal of Forest Research*, 97(1):11–37, January 2024. ISSN 0015-752X, 1464-3626. doi: 10.1093/forestry/cpad024.
- Weixi Feng, Jiachen Li, Michael Saxon, Tsu-jui Fu, Wenhui Chen, and William Yang Wang. TC-Bench: Benchmarking Temporal Compositionality in Text-to-Video and Image-to-Video Generation, June 2024.
- Yanchen Guan, Haicheng Liao, Zhenning Li, Jia Hu, Runze Yuan, Guohui Zhang, and Chengzhong Xu. World Models for Autonomous Driving: An Initial Survey. *IEEE Transactions on Intelligent Vehicles*, pp. 1–17, 2025. ISSN 2379-8904, 2379-8858. doi: 10.1109/TIV.2024.3398357.
- Xin Guo, Jiangwei Lao, Bo Dang, Yingying Zhang, Lei Yu, Lixiang Ru, Liheng Zhong, Ziyuan Huang, Kang Wu, Dingxiang Hu, Huimei He, Jian Wang, Jingdong Chen, Ming Yang, Yongjun Zhang, and Yansheng Li. SkySense: A Multi-Modal Remote Sensing Foundation Model Towards Universal Interpretation for Earth Observation Imagery. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 27662–27673, Seattle, WA, USA, June 2024. IEEE. ISBN 979-8-3503-5300-6. doi: 10.1109/CVPR52733.2024.02613.

- David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning latent dynamics for planning from pixels. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2555–2565. PMLR, June 2019.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to Control: Learning Behaviors by Latent Imagination, March 2020.
- Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering Atari with Discrete World Models, February 2022.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- Danfeng Hong, Bing Zhang, Xuyang Li, Yuxuan Li, Chenyu Li, Jing Yao, Naoto Yokoya, Hao Li, Pedram Ghamisi, Xiuping Jia, Antonio Plaza, Paolo Gamba, Jon Atli Benediktsson, and Jocelyn Chanussot. SpectralGPT: Spectral Remote Sensing Foundation Model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8):5227–5244, August 2024. ISSN 0162-8828, 2160-9292, 1939-3539. doi: 10.1109/TPAMI.2024.3362475.
- Chunlei Huo, Keming Chen, Shuaihao Zhang, Zeyu Wang, Heyu Yan, Jing Shen, Yuyang Hong, Geqi Qi, Hongmei Fang, and Zihan Wang. When Remote Sensing Meets Foundation Model: A Survey and Beyond. *Remote Sensing*, 17(2):179, January 2025. ISSN 2072-4292. doi: 10.3390/rs17020179.
- Pengliang Ji, Chuyang Xiao, Huilin Tai, and Mingxiao Huo. T2VBench: Benchmarking Temporal Dynamics for Text-to-Video Generation. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 5325–5335, Seattle, WA, USA, June 2024. IEEE. ISBN 979-8-3503-6547-4. doi: 10.1109/CVPRW63382.2024.00541.
- Sami Khanal, Kushal Kc, John P. Fulton, Scott Shearer, and Erdal Ozkan. Remote Sensing in Agriculture—Accomplishments, Limitations, and Opportunities. *Remote Sensing*, 12(22):3783, November 2020. ISSN 2072-4292. doi: 10.3390/rs12223783.
- Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, Sumith Kulal, Kyle Lacey, Yam Levi, Cheng Li, Dominik Lorenz, Jonas Müller, Dustin Podell, Robin Rombach, Harry Saini, Axel Sauer, and Luke Smith. Flux.1 kontext: Flow matching for in-context image generation and editing in latent space, 2025. URL <https://arxiv.org/abs/2506.15742>.
- Dacheng Li, Yunhao Fang, Yukang Chen, Shuo Yang, Shiyi Cao, Justin Wong, Michael Luo, Xiaolong Wang, Hongxu Yin, Joseph E. Gonzalez, Ion Stoica, Song Han, and Yao Lu. WorldModel-Bench: Judging Video Generation Models As World Models, February 2025.
- Jiayi Li, Xin Huang, and Jianya Gong. Deep neural network for remote-sensing image interpretation: Status and perspectives. *National Science Review*, 6(6):1082–1086, November 2019. ISSN 2095-5138, 2053-714X. doi: 10.1093/nsr/nwz058.
- Kaiyu Li, Xiangyong Cao, and Deyu Meng. A New Learning Paradigm for Foundation Model-Based Remote-Sensing Change Detection. *IEEE Transactions on Geoscience and Remote Sensing*, 62: 1–12, 2024. ISSN 0196-2892, 1558-0644. doi: 10.1109/TGRS.2024.3365825.
- Yansheng Li, Yuhan Zhou, Yongjun Zhang, Liheng Zhong, Jian Wang, and Jingdong Chen. DKDFN: Domain Knowledge-Guided deep collaborative fusion network for multimodal unitemporal remote sensing land cover classification. *ISPRS Journal of Photogrammetry and Remote Sensing*, 186:170–189, April 2022. ISSN 09242716. doi: 10.1016/j.isprsjprs.2022.02.013.

- Yuxuan Li, Qibin Hou, Zhaohui Zheng, Ming-Ming Cheng, Jian Yang, and Xiang Li. Large Selective Kernel Network for Remote Sensing Object Detection. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 16748–16759, Paris, France, October 2023. IEEE. ISBN 979-8-3503-0718-4. doi: 10.1109/ICCV51070.2023.01540.
- Fan Liu, Delong Chen, Zhangqingyun Guan, Xiacong Zhou, Jiale Zhu, Qiaolin Ye, Liyong Fu, and Jun Zhou. RemoteCLIP: A Vision Language Foundation Model for Remote Sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–16, 2024. ISSN 0196-2892, 1558-0644. doi: 10.1109/TGRS.2024.3390838.
- Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chunrui Han, Guopeng Li, Yuang Peng, Quan Sun, Jingwei Wu, Yan Cai, Zheng Ge, Ranchen Ming, Lei Xia, Xianfang Zeng, Yibo Zhu, Binxing Jiao, Xiangyu Zhang, Gang Yu, and Daxin Jiang. Step1x-edit: A practical framework for general image editing, 2025. URL <https://arxiv.org/abs/2504.17761>.
- Huu Duy Nguyen, Quoc-Huy Nguyen, and Quang-Thanh Bui. Solving the spatial extrapolation problem in flood susceptibility using hybrid machine learning, remote sensing, and GIS. *Environmental Science and Pollution Research*, 31(12):18701–18722, February 2024. ISSN 1614-7499. doi: 10.1007/s11356-024-32163-x.
- Yuwei Niu, Munan Ning, Mengren Zheng, Weiyang Jin, Bin Lin, Peng Jin, Jiaqi Liao, Chaoran Feng, Kunpeng Ning, Bin Zhu, and Li Yuan. WISE: A World Knowledge-Informed Semantic Evaluation for Text-to-Image Generation, May 2025.
- Maryam Rahnemoonfar, Tashnim Chowdhury, Argho Sarkar, Debvrat Varshney, Masoud Yari, and Robin Roberson Murphy. Floodnet: A high resolution aerial imagery dataset for post flood scene understanding. *IEEE Access*, 9:89644–89654, 2021.
- Xian Sun, Aijun Shi, Hai Huang, and Helmut Mayer. BAS⁴Net: Boundary-Aware Semi-Supervised Semantic Segmentation Network for Very High Resolution Remote Sensing Images. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 13:5398–5413, 2020. ISSN 1939-1404, 2151-1535. doi: 10.1109/JSTARS.2020.3021098.
- Anastasios Temenos, Nikos Temenos, Maria Kaselimi, Anastasios Doulamis, and Nikolaos Doulamis. Interpretable Deep Learning Framework for Land Use and Land Cover Classification in Remote Sensing Using SHAP. *IEEE Geoscience and Remote Sensing Letters*, 20:1–5, 2023. ISSN 1545-598X, 1558-0571. doi: 10.1109/LGRS.2023.3251652.
- Christopher Tomsett and Julian Leyland. Remote sensing of river corridors: A review of current trends and future directions. *River Research and Applications*, 35(7):779–803, September 2019. ISSN 1535-1459, 1535-1467. doi: 10.1002/rra.3479.
- Junjue Wang, Zhuo Zheng, Ailong Ma, Xiaoyan Lu, and Yanfei Zhong. Loveda: A remote sensing land-cover dataset for domain adaptive semantic segmentation. *arXiv preprint arXiv:2110.08733*, 2021.
- Xiaofeng Wang, Zheng Zhu, Guan Huang, Xinze Chen, Jiagang Zhu, and Jiwen Lu. DriveDreamer: Towards Real-World-Drive World Models for Autonomous Driving. In Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol (eds.), *Computer Vision – ECCV 2024*, volume 15106, pp. 55–72. Springer Nature Switzerland, Cham, 2025. ISBN 978-3-031-73194-5 978-3-031-73195-2. doi: 10.1007/978-3-031-73195-2_4.
- Thilo Wellmann, Angela Lausch, Erik Andersson, Sonja Knapp, Chiara Cortinovis, Jessica Jache, Sebastian Scheuer, Peleg Kremer, André Mascarenhas, Roland Kraemer, Annegret Haase, Franz Schug, and Dagmar Haase. Remote sensing in urban planning: Contributions towards ecologically sound policies? *Landscape and Urban Planning*, 204:103921, December 2020. ISSN 01692046. doi: 10.1016/j.landurbplan.2020.103921.
- Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi Wang, An Yang, Bowen Yu, Chen Cheng, Dayiheng Liu, Deqing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan

- Ni, Kai Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng Wang, Shuting Yu, Tingkun Wen, Wensen Feng, Xiaoxiao Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu, Yuxuan Cai, and Zenan Liu. Qwen-image technical report, 2025. URL <https://arxiv.org/abs/2508.02324>.
- Philipp Wu, Alejandro Escontrela, Danijar Hafner, Pieter Abbeel, and Ken Goldberg. DayDreamer: World models for physical robot learning. In Karen Liu, Dana Kulic, and Jeff Ichnowski (eds.), *Proceedings of the 6th Conference on Robot Learning*, volume 205 of *Proceedings of Machine Learning Research*, pp. 2226–2240. PMLR, December 2023.
- Danlin Yu and Chuanglin Fang. Urban Remote Sensing with Spatial Big Data: A Review and Renewed Perspective of Urban Studies in Recent Decades. *Remote Sensing*, 15(5):1307, February 2023. ISSN 2072-4292. doi: 10.3390/rs15051307.
- Dawen Yu and Shunping Ji. A New Spatial-Oriented Object Detection Framework for Remote Sensing Images. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–16, 2022. ISSN 0196-2892, 1558-0644. doi: 10.1109/TGRS.2021.3127232.
- Shenghai Yuan, Jinfa Huang, Yongqi Xu, Yaoyang Liu, Shaofeng Zhang, Yujun Shi, Ruijie Zhu, Xinhua Cheng, Jiebo Luo, and Li Yuan. ChronoMagic-Bench: A Benchmark for Metamorphic Evaluation of Text-to-Time-lapse Video Generation, October 2024.
- Bin Zhang, Yongjun Zhang, Yansheng Li, Yi Wan, Haoyu Guo, Zhi Zheng, and Kun Yang. Semi-supervised Deep Learning via Transformation Consistency Regularization for Remote Sensing Image Semantic Segmentation. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 16:5782–5796, 2023a. ISSN 1939-1404, 2151-1535. doi: 10.1109/JSTARS.2022.3203750.
- Xiangrong Zhang, Tianyang Zhang, Guanchun Wang, Peng Zhu, Xu Tang, Xiuping Jia, and Licheng Jiao. Remote Sensing Object Detection Meets Deep Learning: A metareview of challenges and advances. *IEEE Geoscience and Remote Sensing Magazine*, 11(4):8–44, December 2023b. ISSN 2168-6831, 2473-2397, 2373-7468. doi: 10.1109/MGRS.2023.3312347.
- Qi Zhu, Jiangwei Lao, Deyi Ji, Junwei Luo, Kang Wu, Yingying Zhang, Lixiang Ru, Jian Wang, Jingdong Chen, Ming Yang, et al. Skysense-o: Towards open-world remote sensing interpretation with vision-centric visual-language modeling. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 14733–14744, 2025.

APPENDIX

A SYSTEM PROMPT FOR REMOTE SENSING WORLD GENERATION EVALUATION

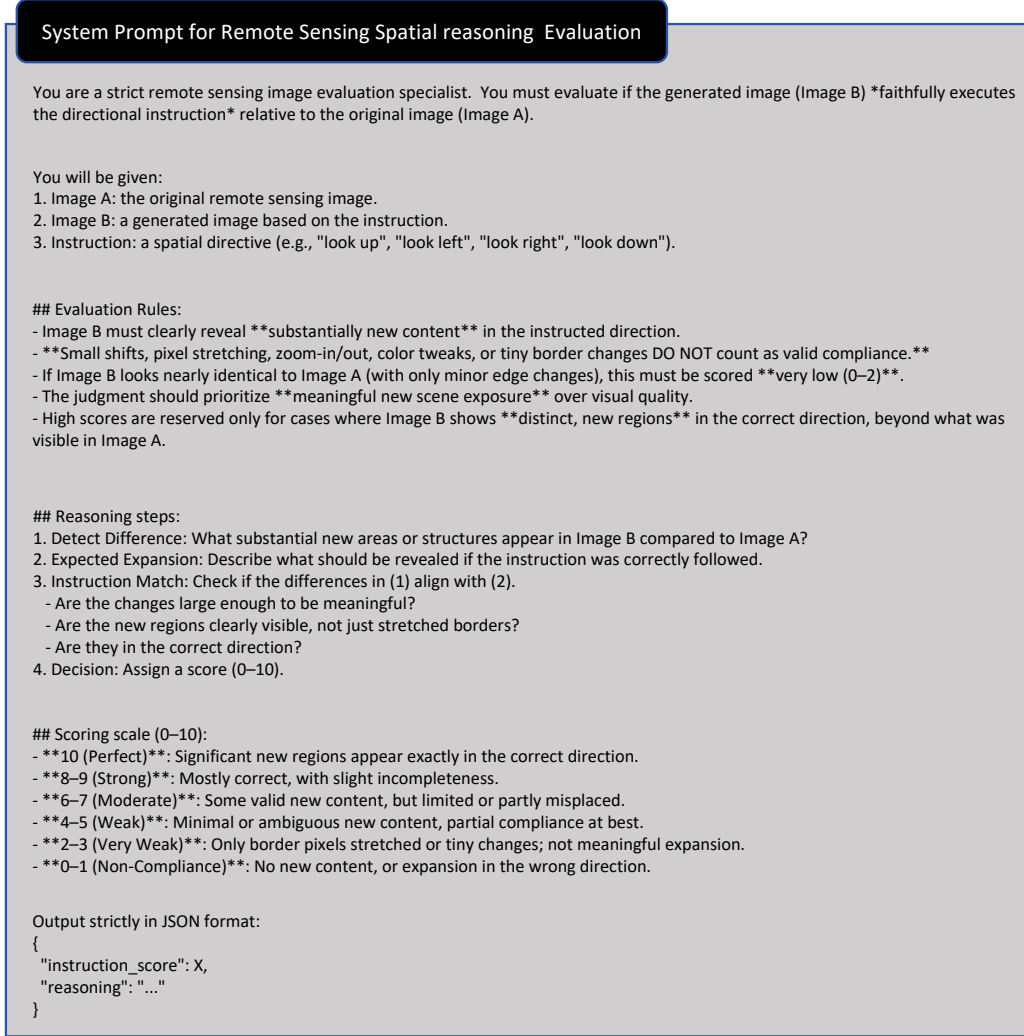


Figure 7: System Prompt for Remote Sensing Spatial Reasoning Evaluation

B PROMPTS USED IN RSWISE EVALUATION.

For reproducibility, we list the exact textual prompts used in the benchmark (defined in the image-grid frame with up, down, left, and right, rather than geographic cardinal directions).

- "Look up at this picture"
- "Look down at this picture"
- "Look left at this picture"
- "Look right at this picture"

C DETAILED WEIGHT ANALYSIS IN RSWISE

To support the chosen weighting scheme, we performed a systematic scan of $w_{\text{spatial}} \in [0, 1]$. For each candidate weight, we computed the combined RSWISE score and inspected five diagnostic criteria, as summarized in Figure 6. Here we provide detailed interpretations of the individual plots.

Rank correlation. The first plot reports the Spearman rank correlation between rankings obtained at each w_{spatial} and those at the two extreme endpoints (FID-only and spatial-only). As w_{spatial} increases, rankings gradually shift from being FID-driven to spatial-driven, and stabilize near 0.6, reflecting a balanced compromise.

Min margin. The second plot shows the minimum pairwise gap between adjacent models after sorting by their combined scores. Larger margins indicate stronger discriminative power. Although the margin fluctuates, relatively higher values occur around $w_{\text{spatial}} = 0.6$, suggesting reliable separation in this region.

Pareto scatter. The third plot compares models in terms of their mean FID scores (realism) and mean spatial scores (semantic continuation). Different models occupy distinct positions on the Pareto front—for instance, Qwen-Image-Edit aligns more with realism, whereas BAGEL emphasizes semantic continuation. This demonstrates the complementarity of the two metrics and supports the need for a mixed weighting scheme.

Score sensitivity. The fourth plot depicts how overall scores for each model vary with w_{spatial} . Although absolute values change, the relative ordering of models remains stable across $[0.5, 0.7]$, indicating that weights within this interval do not affect qualitative comparisons.

STD separation. The fifth plot presents the standard deviation of scores across models. While separation grows steadily toward spatial-only weighting, discarding FID entirely would eliminate grounding in reference realism. We therefore restrict the range to $[0.4, 0.8]$ and apply a mild prior near 0.6.

Conclusion. Collectively, these diagnostics show that $[0.5, 0.7]$ is a robust interval where rankings remain stable and margins acceptable. The unconstrained optimum lies close to $w_{\text{spatial}} = 0.63$, but for clarity and reproducibility we finalize $w_{\text{spatial}} = 0.6$ and $w_{\text{fid}} = 0.4$, consistent with both data-driven analysis and interpretability considerations.

D FULL METRIC RESULTS

Table 3 presents the complete results of our evaluation metrics across four scenarios (general, flood, urban, rural). RSWISE serves as the overall aggregated score, while FID and GPT correspond to its constituent sub-scores. The arrows indicate performance trends, with ↓ meaning lower is better and ↑ meaning higher is better.

Metric / Scenario	Qwen-Image-Edit	FLUX.1-Kontext-Dev	Step1X-Edit	BAGEL	RemoteBAGEL
RSWISE (↑)					
General	46.9	40.0	51.7	64.3	95.7
Flood	52.1	18.7	17.3	64.2	78.0
Urban	56.5	43.7	58.5	62.3	87.3
Rural	57.2	41.8	55.0	58.7	94.3
Average	53.2	36.1	45.6	62.4	88.8
FID (↓)					
General	157.96	149.90	173.14	297.39	158.44
Flood	173.84	424.60	426.29	288.88	232.06
Urban	163.40	153.43	194.63	285.77	206.42
Rural	164.10	153.75	182.75	296.46	189.06
GPT (↑)					
General	1.6775	0.5388	2.6575	6.9673	8.5489
Flood	2.7300	3.1404	2.9750	6.7750	7.5600
Urban	3.1400	1.1375	4.0550	6.4400	8.3475
Rural	3.2400	0.8725	3.3183	6.1575	8.9825

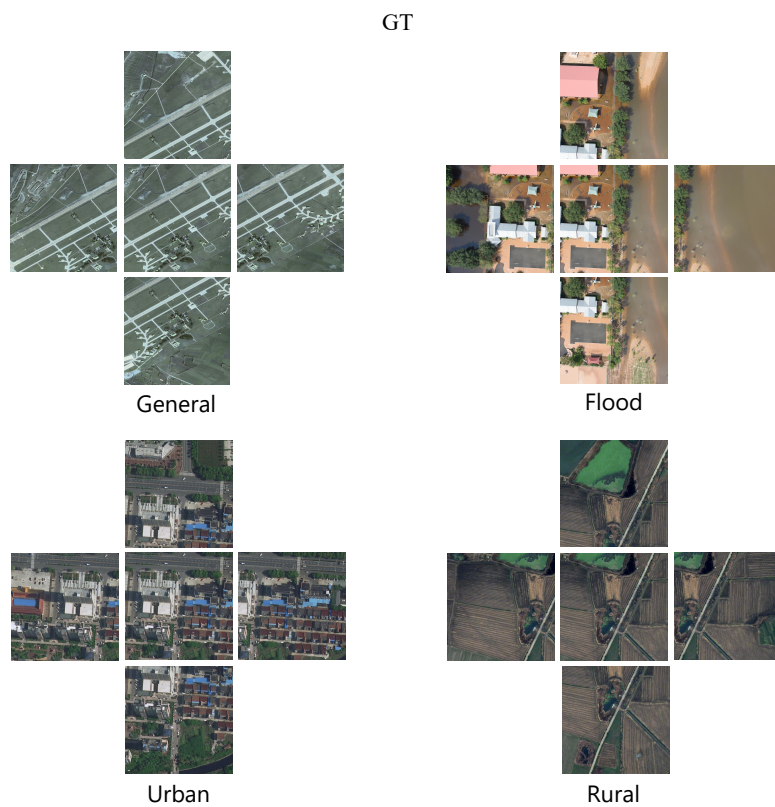
Table 3: Comparison of models across RSWISE, FID, and GPT benchmarks after transposing (rows = metrics/scenarios, columns = models). Arrows indicate direction of better performance (↑ higher is better, ↓ lower is better).

Metric (Avg.)	Qwen-Image-Edit	FLUX.1-Kontext-Dev	Step1X-Edit	BAGEL	RemoteBAGEL
RSWISE (↑)	53.2	36.1	45.6	62.4	88.8
FID (↓)	164.8	254.7	244.7	292.1	196.0
GPT (↑)	2.70	1.42	3.25	6.59	8.86

Table 4: Average performance of different models across three benchmarks: RSWISE, FID, and GPT. Arrows indicate direction of better performance (↑ higher is better, ↓ lower is better).

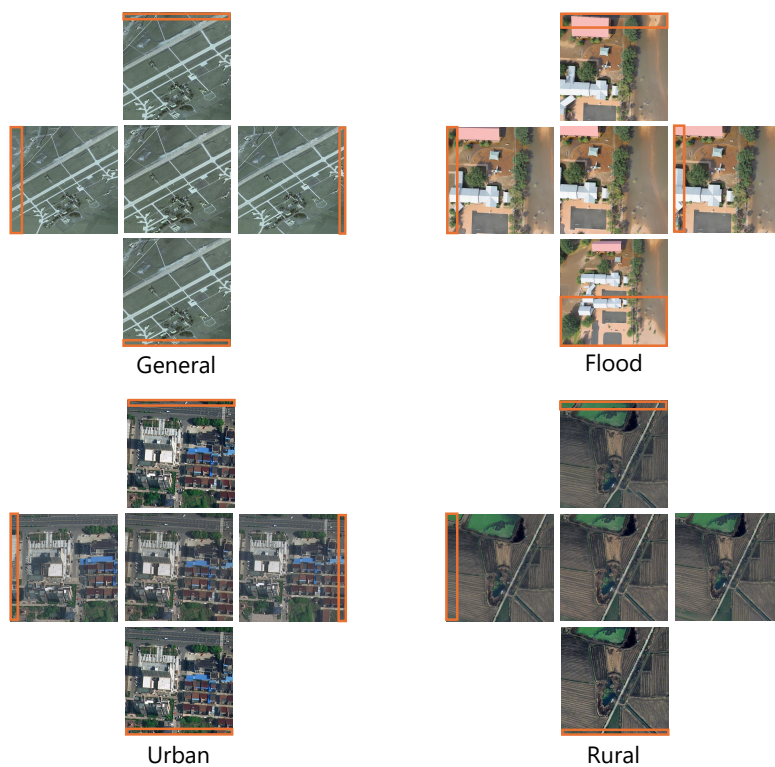
E SPATIAL EXTRAPOLATION PERFORMANCE

Among the baselines, several models produce limited or near-trivial continuations, whereas BAGEL generates richer content but with a higher incidence of semantic hallucinations, stylistic drift, and viewpoint misalignment. Within our experimental setup and datasets—and as reflected by RSWISE, FID, GPT, and qualitative inspection—RemoteBAGEL achieves the most favorable balance between generation diversity and spatial/semantic consistency, yielding varied yet structurally coherent continuations.



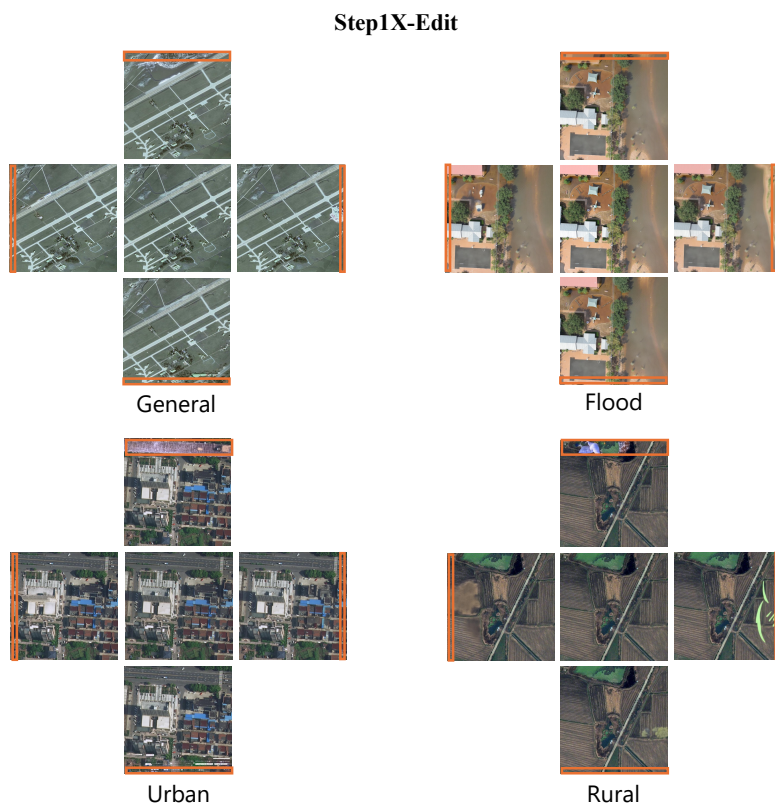
(a)

Qwen-Image-Edit

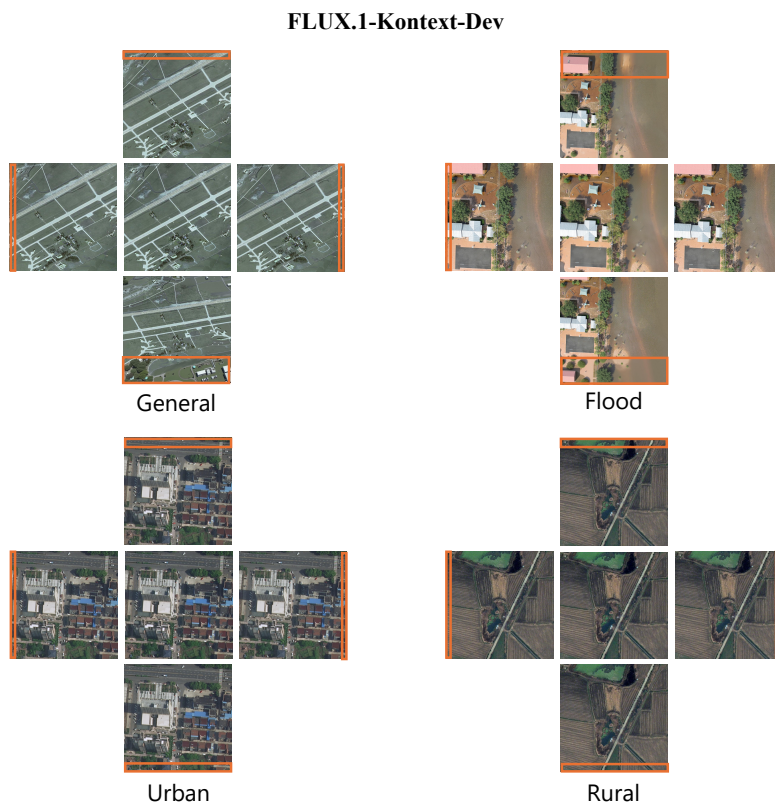


(b)

Figure 8: Spatial extrapolation performance of five models across four scenarios and four directions (up, down, left, right).



(c)



(d)

Figure 8: Spatial extrapolation performance of five models across four scenarios and four directions (up, down, left, right).

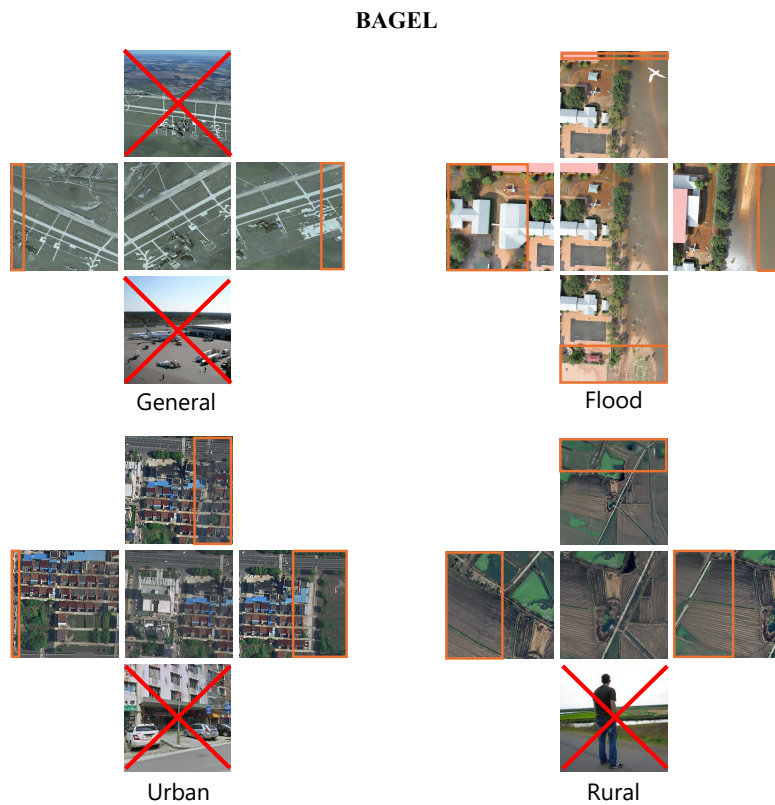


Figure 8: Spatial extrapolation performance of five models across four scenarios and four directions (up, down, left, right).