

Novelty Assessment Report

Paper: Compositional Visual Planning via Inference-Time Diffusion Scaling

PDF URL: <https://openreview.net/pdf?id=EEONns7ae4>

Venue: ICLR 2026 Conference Submission

Year: 2026

Report Generated: 2026-01-05

Abstract

Diffusion models excel at short-horizon robot planning, yet scaling them to long-horizon tasks remains challenging due to computational constraints and limited training data. Existing compositional approaches stitch together short segments by separately denoising each component and averaging overlapping regions. However, this suffers from instability as the factorization assumption breaks down in noisy data space, leading to inconsistent global plans. We propose that the key to stable compositional generation lies in enforcing boundary agreement on the estimated clean data (Tweedie estimates) rather than on noisy intermediate states. Our method formulates long-horizon planning as inference over a chain-structured factor graph of overlapping video chunks, where pretrained short-horizon video diffusion models provide local priors. At inference time, we enforce boundary agreement through a novel combination of synchronous and asynchronous message passing that operates on Tweedie estimates, producing globally consistent guidance without requiring additional training. Our training-free framework demonstrates significant improvements over existing baselines across 100 simulation tasks spanning 4 diverse scenes, effectively generalizing to unseen start-goal combinations that were not present in the original training data. Project website: <https://comp-visual-planning.github.io/>

Disclaimer

This report is **AI-GENERATED** using Large Language Models and WisPaper (a scholar search engine). It analyzes academic papers' tasks and contributions against retrieved prior work. While this system identifies **POTENTIAL** overlaps and novel directions, **ITS COVERAGE IS NOT EXHAUSTIVE AND JUDGMENTS ARE APPROXIMATE**. These results are intended to assist human reviewers and **SHOULD NOT** be relied upon as a definitive verdict on novelty.

Note that some papers exist in multiple, slightly different versions (e.g., with different titles or URLs). The system may retrieve several versions of the same underlying work. The current automated pipeline does not reliably align or distinguish these cases, so human reviewers will need to disambiguate them manually.

If you have any questions, please contact: mingzhang23@m.fudan.edu.cn

Core Task Landscape

This paper addresses: **Long-Horizon Visual Planning Through Compositional Diffusion Models**

A total of **32 papers** were analyzed and organized into a taxonomy with **20 categories**.

Taxonomy Overview

The research landscape has been organized into the following main categories:

- **Trajectory Composition and Stitching Methods**
- **Hierarchical Skill-Based Planning**
- **Constraint-Based and Compositional Planning**
- **Vision-Language Guided Planning**
- **Spatiotemporal and Visuomotor Policy Learning**
- **Foundational Diffusion Planning Frameworks**
- **Specialized Diffusion Applications**

Complete Taxonomy Tree

- Long-Horizon Visual Planning Through Compositional Diffusion Models Survey Taxonomy
- Trajectory Composition and Stitching Methods
 - Overlapping Chunk Composition ★ (2 papers)
 - [0] Compositional Visual Planning via Inference-Time Diffusion Scaling (Anon et al., 2026) [View paper](#)
 - [3] Generative trajectory stitching through diffusion composition (Luo, 2025) [View paper](#)
 - Progressive Trajectory Extension (1 papers)
 - [1] Extendable Planning via Multiscale Diffusion (Chen Chang, 2025) [View paper](#)
- Hierarchical Skill-Based Planning
 - Skill Abstraction and Chaining (2 papers)
 - [2] Generative Skill Chaining: Long-Horizon Skill Planning with Diffusion Models (Mishra, 2024) [View paper](#)
 - [4] SkillDiffuser: Interpretable Hierarchical Planning via Skill Abstractions in Diffusion-Based Task Execution (Zhixuan Liang, 2023) [View paper](#)
 - Coupled Hierarchical Diffusion (1 papers)
 - [14] CHD: Coupled Hierarchical Diffusion for Long-Horizon Tasks (Hao Ce, 2025) [View paper](#)
 - Temporal Logic and Options-Based Planning (2 papers)
 - [9] Diffusion Meets Options: Hierarchical Generative Skill Composition for Temporally-Extended Tasks (Zeyu Feng, 2024) [View paper](#)
 - [10] LTLDoG: Satisfying Temporally-Extended Symbolic Constraints for Safe Diffusion-Based Planning (Zeyu Feng, 2024) [View paper](#)
- Constraint-Based and Compositional Planning
 - Compositional Constraint Satisfaction (2 papers)
 - [16] Compositional Diffusion Models for Powered Descent Trajectory Generation with Flexible Constraints (Julia Briden, 2024) [View paper](#)
 - [21] Compositional Diffusion-Based Continuous Constraint Solvers (Zhutian, 2023) [View paper](#)
 - Scene Graph and Atomic Skill Composition (2 papers)
 - [19] Compose by Focus: Scene Graph-based Atomic Skills (Qi Han, 2025) [View paper](#)
 - [31] Compositional Foundation Models for Hierarchical Planning (Ajay, 2023) [View paper](#)

- Vision-Language Guided Planning
 - Procedure Planning from Visual Observations (2 papers)
 - [17] CLAD: Constrained Latent Action Diffusion for Vision-Language Procedure Planning (Shi Lei, 2025) [View paper](#)
 - [25] Masked Temporal Interpolation Diffusion for Procedure Planning in Instructional Videos (Zhou Yu-fan, 2025) [View paper](#)
 - Vision-Language Navigation (4 papers)
 - [12] Ground Slow, Move Fast: A Dual-System Foundation Model for Generalizable Vision-and-Language Navigation (Meng Wei, 2025) [View paper](#)
 - [18] Trajectory diffusion for objectgoal navigation (Shuqiang Jiang, 2024) [View paper](#)
 - [28] VENTURA: Adapting Image Diffusion Models for Unified Task Conditioned Navigation (Zhang, 2025) [View paper](#)
 - [29] VISTAv2: World Imagination for Indoor Vision-and-Language Navigation (Yanjia Huang, 2025) [View paper](#)
 - Autonomous Driving with Vision-Language Models (2 papers)
 - [8] DiffVLA: Vision-Language Guided Diffusion Planning for Autonomous Driving (Jiang An-qing, 2025) [View paper](#)
 - [13] dVLM-AD: Enhance Diffusion Vision-Language-Model for Driving via Controllable Reasoning (Yingzi Ma, 2025) [View paper](#)
- Spatiotemporal and Visuomotor Policy Learning
 - Spatiotemporal Aware Visuomotor Policies (1 papers)
 - [5] Spatial-Temporal Aware Visuomotor Diffusion Policy Learning (Liu ZhenYang, 2025) [View paper](#)
 - Goal-Conditioned Visual Planning (1 papers)
 - [23] Envision: Embodied Visual Planning via Goal-Imagery Video Diffusion (Yuming Gu, 2025) [View paper](#)
 - Temporal Diffusion Planning (1 papers)
 - [20] Efficient Diffusion Planning with Temporal Diffusion (Jiaming Guo, 2025) [View paper](#)
- Foundational Diffusion Planning Frameworks
 - Trajectory Optimization as Modeling (1 papers)
 - [7] Planning with Diffusion for Flexible Behavior Synthesis (Janner, 2022) [View paper](#)
 - Diffusion Models in Robotics Survey (1 papers)
 - [11] Diffusion Models in Robotics: A Survey (X Liu, 2025) [View paper](#)
- Specialized Diffusion Applications
 - Long-Term Motion and Action Generation (2 papers)
 - [6] Gated temporal diffusion for stochastic long-term dense anticipation (Olga Zatsarynna, 2024) [View paper](#)
 - [22] Synthesizing Long-Term Human Motions with Diffusion Models via Coherent Sampling (Yang Zhao, 2023) [View paper](#)
 - Human-Robot Collaboration and Co-Policy Learning (1 papers)
 - [15] Long-Horizon Prediction for Human-Robot Collaboration (Ng, 2023) [View paper](#)
 - One-Shot Compositional Subgoal Learning (1 papers)
 - [30] Generalizing to New Tasks via One-Shot Compositional Subgoals (Bian Xihan, 2022) [View paper](#)
 - Diffusion-Based Vision-Language Models (1 papers)
 - [24] Dream-VL & Dream-VLA: Open Vision-Language and Vision-Language-Action Models with Diffusion Language Model Backbone (Jiacheng Ye, 2025) [View paper](#)
 - Compositional Visual and Text Generation (3 papers)
 - [26] Compositional Visual Generation With Enhanced Language Guidance (Feng, 2025) [View paper](#)
 - [27] Loom: Diffusion-Transformer for Interleaved Generation (Mingcheng Ye, 2025) [View paper](#)
 - [32] ComposeAnything: Composite Object Priors for Text-to-Image Diffusion Models (Z Khan, n.d.) [View paper](#)

Narrative

Core task: long-horizon visual planning through compositional diffusion models. The field addresses how to generate extended action sequences or visual trajectories by leveraging diffusion-based generative models that can compose or stitch together shorter segments. The taxonomy reveals several complementary directions: Trajectory Composition and Stitching Methods focus on connecting overlapping or adjacent trajectory chunks to extend planning horizons, as seen in works like Generative Trajectory Stitching[3] and Generative Skill Chaining[2]. Hierarchical Skill-Based Planning decomposes tasks into reusable primitives, exemplified by SkillDiffuser[4] and related approaches. Constraint-Based and Compositional Planning emphasizes satisfying logical or geometric constraints during generation, while Vision-Language Guided Planning integrates natural language instructions to steer diffusion processes. Spatiotemporal and Visuomotor Policy Learning targets direct sensorimotor control with temporal coherence, and Foundational Diffusion Planning Frameworks provide core algorithmic innovations such as Planning with Diffusion[7]. Specialized Diffusion Applications explore domain-specific uses ranging from robotics to autonomous driving.

Within this landscape, a particularly active line of work centers on trajectory stitching and chunk composition, where the challenge is to seamlessly merge locally generated segments into coherent long-horizon plans. Compositional Visual Planning[0] sits squarely in this branch, employing overlapping chunk composition to extend planning reach beyond what single-shot diffusion models can achieve. Its approach closely aligns with Generative Trajectory Stitching[3], which also addresses how to blend trajectory pieces, though the two may differ in their blending mechanisms or the granularity of overlap. Meanwhile, hierarchical methods like SkillDiffuser[4] offer an alternative by learning discrete skill libraries, trading off the flexibility of continuous stitching for the interpretability of modular primitives. Across these branches, open questions remain about how to balance computational efficiency, sample quality, and the ability to handle diverse constraints—issues that Compositional Visual Planning[0] and its neighbors continue to explore through different compositional strategies.

Related Works in Same Category

The following **1 sibling papers** share the same taxonomy leaf node with the original paper:

1. Generative trajectory stitching through diffusion composition

Authors: Luo, Yunhao, Mishra, Utkarsh A., Du, et al. (8 authors total) | **Year/Venue:** 2025 | **URL:** [View paper](#)

Abstract

Effective trajectory stitching for long-horizon planning is a significant challenge in robotic decision-making. While diffusion models have shown promise in planning, they are limited to solving tasks similar to those seen in their training data. We propose CompDiffuser, a novel generative approach that can solve new tasks by learning to compositionally stitch together shorter trajectory chunks from previously seen tasks. Our key insight is modeling the trajectory distribution by subdividing it ...

Relationship Analysis

Both papers belong to the Overlapping Chunk Composition category, addressing long-horizon visual planning by subdividing trajectories into overlapping segments and learning conditional relationships for compositional generation. The original paper focuses on inference-

time message passing (synchronous and asynchronous) on Tweedie estimates to enforce boundary agreement in a chain-structured factor graph, while the candidate paper (CompDiffuser) trains a single bidirectional diffusion model that conditions on noisy samples of neighboring chunks during the denoising process, enabling both parallel and autoregressive sampling schemes. The key difference lies in the compositional mechanism: the original paper operates on clean Tweedie estimates with training-free guidance, whereas the candidate conditions directly on noisy neighboring chunks during training and sampling.

Contributions Analysis

Overall novelty summary. The paper proposes a training-free compositional framework for long-horizon visual planning that enforces boundary agreement on Tweedie estimates rather than noisy intermediate states. It resides in the 'Overlapping Chunk Composition' leaf under 'Trajectory Composition and Stitching Methods', which contains only two papers total. This places the work in a relatively sparse research direction within the broader taxonomy of 32 papers across multiple branches. The sibling paper in this leaf, Generative Trajectory Stitching, also addresses overlapping chunk composition, suggesting that this specific approach to long-horizon planning is an emerging but not yet crowded area.

The taxonomy reveals that neighboring leaves include 'Progressive Trajectory Extension' (one paper) and broader sibling branches like 'Hierarchical Skill-Based Planning' (six papers across three sub-categories) and 'Constraint-Based and Compositional Planning' (four papers). The paper's focus on factor graph inference over video chunks distinguishes it from hierarchical skill decomposition methods, which learn discrete primitives, and from constraint satisfaction approaches that compose energies. The scope note for this leaf explicitly excludes progressive extension without overlap and multiscale hierarchical methods, clarifying that the paper's overlapping chunk strategy occupies a distinct methodological niche within trajectory composition.

Among 24 candidates examined across three contributions, the analysis found limited prior work overlap. The core contribution of boundary agreement on Tweedie estimates examined four candidates with zero refutable matches. The message passing mechanism examined ten candidates, also with zero refutable matches, suggesting novelty in the inference procedure. However, the compositional planning benchmark contribution examined ten candidates and found two refutable matches, indicating that evaluation frameworks for compositional generalization may have more substantial prior work. The limited search scope means these findings reflect top-K semantic matches rather than exhaustive coverage.

Based on the limited literature search of 24 candidates, the work appears to introduce novel inference mechanisms within a sparse research direction. The taxonomy structure shows that overlapping chunk composition itself is an emerging area with few direct comparisons. The two refutable matches for the benchmark contribution suggest that evaluation methodologies may be less novel than the core algorithmic approach, though the restricted search scope prevents definitive conclusions about the broader landscape.

This paper presents **3 main contributions**, each analyzed against relevant prior work:

Contribution 1: Compositional visual planning via boundary agreement on Tweedie estimates

Description: The authors formulate long-horizon planning as inference over a chain-structured factor graph of overlapping video chunks. Instead of enforcing consistency on noisy diffusion states (as in prior work), they enforce boundary agreement on Tweedie estimates (estimated clean data), addressing the core limitation that factorization assumptions break down during diffusion sampling.

This contribution was assessed against **4 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. TweedieMix: Improving Multi-Concept Fusion for Diffusion-based Image/Video Generation

URL: [View paper](#)

Brief Assessment

TweedieMix[53] focuses on multi-concept fusion for image/video generation by mixing denoised outputs in Tweedie space for appearance blending. The original paper addresses long-horizon planning via boundary agreement constraints in factor graphs for temporal consistency. These are fundamentally different applications of Tweedie estimates.

2. Improved Sampling Of Diffusion Models In Fluid Dynamics With Tweedie's Formula

URL: [View paper](#)

Brief Assessment

Improved Sampling Fluid Dynamics[52] focuses on accelerating diffusion model sampling for fluid dynamics simulations through truncated sampling and iterative refinement, not on compositional generation or boundary agreement for long-horizon planning tasks.

3. Motion Composition and Interpolation Using Diffusion Models

URL: [View paper](#)

Brief Assessment

Motion Composition Interpolation[55] focuses on blending robot motion primitives in configuration space using weighted score averaging for interpolation, not on enforcing boundary agreement constraints for long-horizon planning via factor graphs.

4. Compositional simulation-based inference for time series

URL: [View paper](#)

Brief Assessment

Compositional Simulation Inference[54] addresses simulation-based inference for time series using Markovian simulators, not visual planning with diffusion models. The candidate focuses on parameter inference from temporal data through local score estimation and composition, while the original paper tackles long-horizon robot planning through video generation with boundary agreement on Tweedie estimates in diffusion models.

Contribution 2: Joint synchronous and asynchronous message passing on denoised variables

Description: The authors introduce two complementary message-passing mechanisms that operate on Tweedie estimates: a synchronous scheme treating the chain as a Gaussian linear system with parallel updates, and an asynchronous scheme using one-sided stop-gradient targets for faster convergence. These are integrated into a training-free DDIM sampler via diffusion-sphere guidance.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. Partially Conditioned Patch Parallelism for Accelerated Diffusion Model Inference

URL: [View paper](#)

Brief Assessment

Patch Parallelism Acceleration[38] focuses on accelerating diffusion model inference through patch-based parallelism with asynchronous communication between neighboring patches, not on message passing mechanisms for compositional planning with Tweedie estimates in factor graphs.

2. Enhancing Approximate Message Passing via Diffusion Models Towards On-Device Intelligence

URL: [View paper](#)

Brief Assessment

Approximate Message Passing[39] focuses on signal recovery and noise estimation in compressed sensing contexts, not on compositional planning with boundary agreement constraints in diffusion models.

3. Deep networks as denoising algorithms: Sample-efficient learning of diffusion models in high-dimensional graphical models

URL: [View paper](#)

Brief Assessment

Deep Networks Denoising[33] focuses on score function approximation in graphical models (Ising models, sparse coding) using variational inference, not on message passing mechanisms for compositional diffusion planning in visual domains.

4. Using Powerful Prior Knowledge of Diffusion Model in Deep Unfolding Networks for Image Compressive Sensing

URL: [View paper](#)

Brief Assessment

Prior Knowledge Compressive Sensing[37] focuses on image compressive sensing using diffusion message passing (DMP) for reconstruction from compressed measurements, not on compositional planning with boundary agreement constraints in factor graphs.

5. DG-RainDiff: Depth-Guided Dynamic Message Passing Diffusion Model for Mixture of Rain Removal

URL: [View paper](#)

Brief Assessment

DG-RainDiff[42] focuses on image deraining with a dynamic message passing module for spatial feature extraction, not on compositional planning with synchronous/asynchronous message passing schemes operating on Tweedie estimates for boundary agreement in factor graphs.

6. Your diffusion model is secretly a noise classifier and benefits from contrastive training

URL: [View paper](#)

Brief Assessment

Noise Classifier Contrastive[41] focuses on improving diffusion model training through contrastive loss for noise classification, not on message passing mechanisms for compositional planning. The candidate addresses denoiser robustness in out-of-distribution regions during sampling, while the original contribution proposes specific synchronous/asynchronous message passing schemes on Tweedie estimates for long-horizon visual planning.

7. Asyncdiff: Parallelizing diffusion models by asynchronous denoising

URL: [View paper](#)

Brief Assessment

Asyncdiff[36] focuses on parallelizing diffusion models through asynchronous denoising across model components distributed on different devices, not on message passing between denoised variables for compositional planning. The asynchronous mechanism in Asyncdiff[36] refers to breaking sequential dependencies in the denoising model itself, whereas the original paper's contribution addresses message passing on Tweedie estimates for boundary agreement in factor graphs.

8. CL-DiffPhyCon: Closed-loop Diffusion Control of Complex Physical Systems

URL: [View paper](#)

Brief Assessment

CL-DiffPhyCon[40] focuses on physical systems control using asynchronous denoising across physical time steps, not on compositional visual planning with boundary agreement on Tweedie estimates for video generation tasks.

9. SAFedHDM: Semi-asynchronous federated learning with highlight diffusion model for medical image segmentation

URL: [View paper](#)

Brief Assessment

SAFedHDM[35] addresses federated learning for medical image segmentation with semi-asynchronous communication protocols, not message passing on denoised variables in diffusion sampling for compositional planning.

10. SCARefusion: Side channel analysis data restoration with diffusion model

URL: [View paper](#)

Brief Assessment

SCARefusion[34] focuses on side channel analysis data restoration using diffusion models for cryptographic security applications, not on compositional planning or message passing mechanisms for enforcing boundary agreement in factor graphs.

Contribution 3: Compositional planning benchmark for evaluating generalization to unseen start-goal combinations

Description: The authors develop a benchmark for compositional planning in robotic manipulation where training data contains only N start-goal pairs, but evaluation includes N·N·N unseen combinations. This tests whether planners can generalize by composing fragments from the training distribution to solve novel tasks.

This contribution was assessed against **10 related papers** from the literature. Papers with potential prior art are analyzed in detail with textual evidence; others receive brief assessments.

1. A Benchmark for Compositional Visual Reasoning

URL: [View paper](#)

Brief Assessment

Compositional Visual Reasoning Benchmark[46] focuses on visual reasoning tasks with abstract relations and object attributes (shape, color, size), not robotic manipulation planning with start-goal state combinations. The benchmark evaluates compositional visual reasoning abilities rather than planning generalization.

2. Exedec: Execution decomposition for compositional generalization in neural program synthesis

URL: [View paper](#)

Brief Assessment

Exedec[44] focuses on compositional generalization in neural program synthesis (string/list manipulation tasks), not robotic manipulation planning with visual observations and start-goal image pairs.

3. Learning from Less: Guiding Deep Reinforcement Learning with Differentiable Symbolic Planning

URL: [View paper](#)

Brief Assessment

Learning from Less[51] focuses on symbolic planning with policy primitives in grid-based environments, not on visual planning benchmarks. The candidate does not present a compositional planning benchmark structure similar to the original paper's N·N-N evaluation framework.

4. Generative trajectory stitching through diffusion composition

URL: [View paper](#)

Prior Art Analysis

Generative Trajectory Stitching[3] demonstrates that a similar benchmark for compositional planning with unseen start-goal combinations was already established. The candidate paper explicitly describes training on N start-goal pairs and evaluating on N·N-N unseen combinations in multiple environments (pointmaze, antmaze, humanoidmaze, antsoccer). This benchmark design predates the original paper's contribution and serves the same purpose of testing whether planners can generalize by composing fragments from the training distribution to solve novel tasks.

Evidence

Evidence 1 - **Rationale:** This confirms that Generative Trajectory Stitching[3] established benchmark tasks for compositional planning across multiple environments and difficulty levels. - **Original:** benchmark for compositional planning in robotic manipulation - **Candidate:** we conduct experiments on benchmark tasks of various difficulties, covering different environment sizes, agent state dimension, trajectory types, training data quality, and show that compdiffuser significantly outperforms existing methods.

5. Language model agents suffer from compositional generalization in web automation

URL: [View paper](#)

Brief Assessment

[Final Audit Failure] The model insisted on a refutation claim but failed to provide verifiable evidence after multiple retries. Marked as cannot_refute for safety. Please manually verify the candidate text.

6. Compositional generalization via neural-symbolic stack machines

URL: [View paper](#)

Brief Assessment

Neural Symbolic Stack Machines[47] focuses on compositional generalization in language-driven navigation and grammar parsing tasks, not on robotic manipulation planning benchmarks with start-goal combinations.

7. Environment generation for zero-shot compositional reinforcement learning

URL: [View paper](#)

Prior Art Analysis

Environment Generation Zero Shot[50] demonstrates prior work on compositional planning benchmarks that evaluate generalization to unseen start-goal combinations. The candidate paper explicitly describes a benchmark where training data contains only N start-goal pairs while evaluation includes N·N-N unseen combinations, directly matching the original paper's contribution. Both papers test whether planners can generalize by composing fragments from the training distribution to solve novel tasks, with the candidate paper providing this framework in 2021, predating the original submission to ICLR 2026.

Evidence

Evidence 1 - **Rationale:** Both papers address the same core problem: evaluating whether agents can generalize to unseen task combinations by composing learned fragments. The candidate explicitly frames this as compositional generalization where agents must combine sub-tasks in novel ways, matching the original's goal of testing generalization to 'cross-region plans unseen in the dataset but composable from its fragments.' - **Original:** for example, in the tool-use setting (figure 6), demonstrations cover only the blue or green regions; the planner must generalize across them to form cross-region plans unseen in the dataset but composable from its fragments. - **Candidate:** many real-world problems are compositional - solving them requires completing interdependent sub-tasks, either in series or in parallel, that can be represented as a dependency graph. deep reinforcement learning (rl) agents often struggle to learn such complex tasks due to the long time horizons and...

Evidence 2 - **Rationale:** Both papers emphasize zero-shot generalization to unseen task combinations at test time. The candidate's framework explicitly evaluates 'zero-shot' generalization to 'unseen tasks at test-time,' which directly corresponds to the original's evaluation of 'unseen start-goal combinations that were not present in the original training data.' - **Original:** we address this type of generation by learning from short demonstration chunks randomly taken from long-horizon tasks and compositionally generates multiple chunks at inference time to construct the final plan. - **Candidate:** this automatic curriculum not only enables the agent to learn more complex tasks than it could have otherwise, but also selects tasks where the agent's performance is weak, enhancing its robustness and ability to generalize zero-shot to unseen tasks at test-time.

8. What Do You Need for Diverse Trajectory Stitching in Diffusion Planning?

URL: [View paper](#)

Brief Assessment

[Final Audit Failure] The model insisted on a refutation claim but failed to provide verifiable evidence after multiple retries. Marked as cannot_refute for safety. Please manually verify the candidate text.

9. Imagine the Unseen World: A Benchmark for Systematic Generalization in Visual World Models

URL: [View paper](#)

Brief Assessment

Imagine Unseen World[48] focuses on systematic visual imagination in one-step image-to-image transformations with object-centric factors, not on compositional planning for robotic manipulation with start-goal pairs in action sequences.

10. AGQA: A Benchmark for Compositional Spatio-Temporal Reasoning

URL: [View paper](#)

Brief Assessment

AGQA[43] focuses on compositional spatio-temporal reasoning for video question answering, not robotic manipulation planning. The benchmark tests visual reasoning over temporal actions and spatial relationships in videos, which is fundamentally different from the original paper's compositional planning framework for robotic manipulation with start-goal combinations.

Appendix: Text Similarity Detection

No high-similarity text segments were detected across any compared papers.

References

- [0] Compositional Visual Planning via Inference-Time Diffusion Scaling [View paper](#)
- [1] Extendable Planning via Multiscale Diffusion [View paper](#)
- [2] Generative Skill Chaining: Long-Horizon Skill Planning with Diffusion Models [View paper](#)
- [3] Generative trajectory stitching through diffusion composition [View paper](#)
- [4] SkillDiffuser: Interpretable Hierarchical Planning via Skill Abstractions in Diffusion-Based Task Execution [View paper](#)
- [5] Spatial-Temporal Aware Visuomotor Diffusion Policy Learning [View paper](#)
- [6] Gated temporal diffusion for stochastic long-term dense anticipation [View paper](#)
- [7] Planning with Diffusion for Flexible Behavior Synthesis [View paper](#)
- [8] DiffVLA: Vision-Language Guided Diffusion Planning for Autonomous Driving [View paper](#)
- [9] Diffusion Meets Options: Hierarchical Generative Skill Composition for Temporally-Extended Tasks [View paper](#)
- [10] LTLDoG: Satisfying Temporally-Extended Symbolic Constraints for Safe Diffusion-Based Planning [View paper](#)
- [11] Diffusion Models in Robotics: A Survey [View paper](#)
- [12] Ground Slow, Move Fast: A Dual-System Foundation Model for Generalizable Vision-and-Language Navigation [View paper](#)
- [13] dVLM-AD: Enhance Diffusion Vision-Language-Model for Driving via Controllable Reasoning [View paper](#)
- [14] CHD: Coupled Hierarchical Diffusion for Long-Horizon Tasks [View paper](#)
- [15] Long-Horizon Prediction for Human-Robot Collaboration [View paper](#)
- [16] Compositional Diffusion Models for Powered Descent Trajectory Generation with Flexible Constraints [View paper](#)
- [17] CLAD: Constrained Latent Action Diffusion for Vision-Language Procedure Planning [View paper](#)
- [18] Trajectory diffusion for objectgoal navigation [View paper](#)
- [19] Compose by Focus: Scene Graph-based Atomic Skills [View paper](#)
- [20] Efficient Diffusion Planning with Temporal Diffusion [View paper](#)
- [21] Compositional Diffusion-Based Continuous Constraint Solvers [View paper](#)
- [22] Synthesizing Long-Term Human Motions with Diffusion Models via Coherent Sampling [View paper](#)
- [23] Envision: Embodied Visual Planning via Goal-Imagery Video Diffusion [View paper](#)
- [24] Dream-VL & Dream-VLA: Open Vision-Language and Vision-Language-Action Models with Diffusion Language Model Backbone [View paper](#)
- [25] Masked Temporal Interpolation Diffusion for Procedure Planning in Instructional Videos [View paper](#)
- [26] Compositional Visual Generation With Enhanced Language Guidance [View paper](#)
- [27] Loom: Diffusion-Transformer for Interleaved Generation [View paper](#)
- [28] VENTURA: Adapting Image Diffusion Models for Unified Task Conditioned Navigation [View paper](#)
- [29] VISTAv2: World Imagination for Indoor Vision-and-Language Navigation [View paper](#)
- [30] Generalizing to New Tasks via One-Shot Compositional Subgoals [View paper](#)
- [31] Compositional Foundation Models for Hierarchical Planning [View paper](#)
- [32] ComposeAnything: Composite Object Priors for Text-to-Image Diffusion Models [View paper](#)
- [33] Deep networks as denoising algorithms: Sample-efficient learning of diffusion models in high-dimensional graphical models [View paper](#)
- [34] SCARefusion: Side channel analysis data restoration with diffusion model [View paper](#)
- [35] SAFedHDM: Semi-asynchronous federated learning with highlight diffusion model for medical image segmentation [View paper](#)
- [36] Asyncdiff: Parallelizing diffusion models by asynchronous denoising [View paper](#)
- [37] Using Powerful Prior Knowledge of Diffusion Model in Deep Unfolding Networks for Image Compressive Sensing [View paper](#)
- [38] Partially Conditioned Patch Parallelism for Accelerated Diffusion Model Inference [View paper](#)
- [39] Enhancing Approximate Message Passing via Diffusion Models Towards On-Device Intelligence [View paper](#)
- [40] CL-DiffPhyCon: Closed-loop Diffusion Control of Complex Physical Systems [View paper](#)
- [41] Your diffusion model is secretly a noise classifier and benefits from contrastive training [View paper](#)
- [42] DG-RainDiff: Depth-Guided Dynamic Message Passing Diffusion Model for Mixture of Rain Removal [View paper](#)
- [43] AGQA: A Benchmark for Compositional Spatio-Temporal Reasoning [View paper](#)
- [44] Exedec: Execution decomposition for compositional generalization in neural program synthesis [View paper](#)
- [45] Language model agents suffer from compositional generalization in web automation [View paper](#)
- [46] A Benchmark for Compositional Visual Reasoning [View paper](#)
- [47] Compositional generalization via neural-symbolic stack machines [View paper](#)
- [48] Imagine the Unseen World: A Benchmark for Systematic Generalization in Visual World Models [View paper](#)
- [49] What Do You Need for Diverse Trajectory Stitching in Diffusion Planning? [View paper](#)
- [50] Environment generation for zero-shot compositional reinforcement learning [View paper](#)
- [51] Learning from Less: Guiding Deep Reinforcement Learning with Differentiable Symbolic Planning [View paper](#)
- [52] Improved Sampling Of Diffusion Models In Fluid Dynamics With Tweedie's Formula [View paper](#)
- [53] TweedieMix: Improving Multi-Concept Fusion for Diffusion-based Image/Video Generation [View paper](#)
- [54] Compositional simulation-based inference for time series [View paper](#)
- [55] Motion Composition and Interpolation Using Diffusion Models [View paper](#)