

# [TINY PAPER] TRAINING-FREE CONSTRUCTION OF EXECUTABLE 3D WORLDS FROM NARRATIVE TEXT

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

A prominent trend in recent world generation and world model research emphasizes foundation-scale approaches, particularly diffusion-based architectures trained on large video and multimodal datasets, often requiring significant computational resources. While effective, this paradigm implicitly assumes access to infrastructure that is unavailable to many researchers and practitioners. In this work, we explore an alternative perspective on world model construction under strict compute constraints. We present a modular, training-free framework that leverages existing multimodal large language models (MLLMs) and open-source text-to-3D asset generators through lightweight API calls to construct story-driven, navigable 3D worlds. Rather than learning world dynamics end-to-end, our system extracts structured semantic representations from narrative text and deterministically compiles them into spatial layouts, connectivity graphs, and executable environments. We demonstrate that coherent, traversable worlds can be generated on commodity hardware, suggesting that world model research can advance not only through scaling compute, but also through structural abstraction, compositional design, and systems-level reasoning.

## 1 INTRODUCTION

Recent advances in world model research have followed a clear generative trajectory, progressing from high-fidelity static image synthesis to temporally coherent video generation, and more recently toward real-time and interactive environments. Latent diffusion models first demonstrated that complex visual distributions could be efficiently modeled in compressed latent spaces (Rombach et al., 2022), establishing diffusion as a scalable paradigm for image-based world generation. This foundation was subsequently extended to the temporal domain, enabling video diffusion models capable of modeling spatiotemporal dynamics and long-horizon visual consistency (Blattmann et al., 2023). Building on these developments, recent systems have begun to incorporate agent conditioning and interaction, positioning diffusion models as general-purpose simulators that respond dynamically to user actions (Bruce et al., 2024). Together, these works reflect a broader shift toward world models that aim to support continuous, interactive, and embodied experience.

Despite their impressive capabilities, this progression has been accompanied by rapidly increasing computational demands. Recent studies have shown that diffusion-based architectures exhibit clear power-law scaling behavior with respect to compute budget, with improvements in generation quality requiring joint increases in model size, data, and total training FLOPs (Liang et al., 2024). In practice, diffusion-based world models for video and interaction incur substantial training and inference costs, particularly as resolution, temporal horizon, and agent conditioning increase (Blattmann et al., 2023; Bruce et al., 2024). As a result, many such systems remain accessible primarily to well-resourced research labs, limiting their practicality for individual researchers, educators, and practitioners.

In this work, we propose an alternative approach to world model construction that prioritizes explicit structure, determinism, and accessibility under strict compute constraints. Rather than learning world dynamics end-to-end, our method constructs story-driven, navigable 3D worlds through a modular, training-free pipeline that combines pretrained multimodal language models accessed via lightweight API calls with algorithmic world assembly. This perspective complements diffusion-

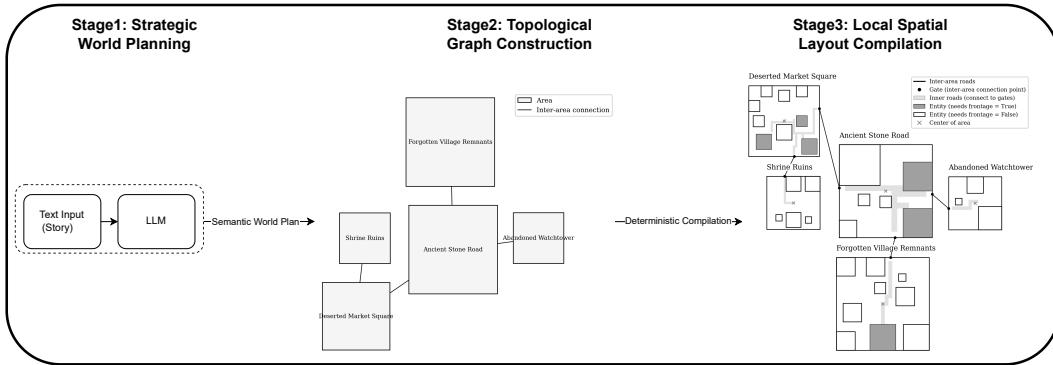


Figure 1: Overview of the training-free world construction pipeline. A natural-language narrative is parsed by a large language model into a semantic world plan that specifies regions and their associated entities (Stage 1). This plan is used to construct an abstract topological world graph defining relative placement and traversable connections between regions (Stage 2). The graph is then deterministically compiled into collision-free, tile-based spatial layouts, with entities grounded within each area and explicit traversal gates and routed road networks (Stage 3).

based world models by opening a distinct design space centered on structured abstraction and compositional reasoning, rather than large-scale training and specialized hardware.

## 2 METHODOLOGY

Our approach constructs a 3D world through five sequential stages, transforming a natural-language narrative into an executable environment instantiated in (Godot Engine Project, 2024). A large language model (GPT-4o (OpenAI, 2024)) is used to extract areas, entities, and relational constraints from text, while all spatial layout, connectivity, and world assembly are handled algorithmically and deterministically.

### 2.1 STAGE 1: STRATEGIC WORLD PLANNING

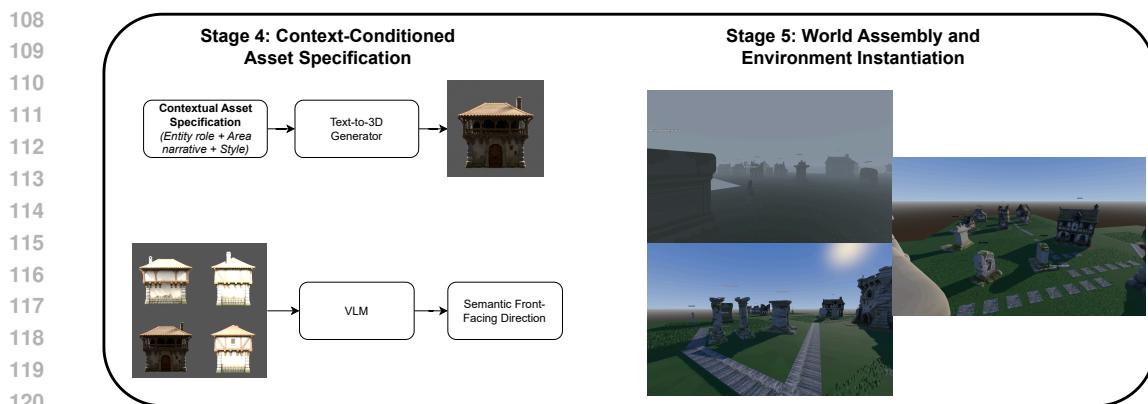
The process begins with a natural-language story as input, which is transformed by a large language model into a structured world plan consisting of semantic areas and required entities. Each area is annotated with a localized narrative description and symbolic scale information, providing high-level constraints on the intended spatial extent and content of the region.

Crucially, the language model is not used to predict geometric coordinates. Instead, spatial grounding is deferred to later stages, ensuring that semantic planning remains flexible while all geometric structure is derived deterministically. This separation avoids unreliable geometric inference from language models and enables reproducible world construction.

### 2.2 STAGE 2: TOPOLOGICAL GRAPH CONSTRUCTION

Given the semantic areas identified in Stage 2.1, we construct a topological world graph that defines relative placement and traversable connectivity between regions. Areas are arranged using symbolic spatial relationships, while explicit inter-area connections specify which regions are directly traversable.

This representation captures global structure without committing to exact geometry. All symbolic placements and connections are grounded into concrete spatial layouts during subsequent deterministic compilation, ensuring both coherent world organization and reproducible connectivity.



121  
122  
123  
124  
125  
126  
127  
128

Figure 2: Context-conditioned asset grounding and world assembly. Given the resolved tile-based world layout, context-conditioned 3D assets are generated for each entity by conditioning text-to-3D prompts on the entity’s semantic role, local area narrative, and global style constraints (Stage 4). Because text-to-3D models produce meshes with arbitrary orientation, canonical front, back, left, and right views are rendered and a vision-language model is used to identify the semantic front-facing direction when required. The resolved tile-based layout and oriented assets are then deterministically assembled into an executable environment instantiated in the Godot engine, producing a coherent and traversable 3D world (Stage 5).

### 131 2.3 STAGE 3: LOCAL SPATIAL LAYOUT COMPILED

132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

For each area in the world graph (Stage 2.2), we deterministically ground entities, traversal interfaces, and road networks onto a discrete two-dimensional tile grid. Each area is assigned a fixed grid resolution based on its symbolic scale and compiled independently to preserve modularity.

Entities and gates are placed under collision and boundary constraints, and a connected road network is generated to ensure full traversability between regions and landmarks. All geometric decisions are resolved algorithmically, yielding a collision-free and reproducible spatial layout that fully specifies world geometry.

### 2.4 STAGE 4: CONTEXT-CONDITIONED ASSET SPECIFICATION

Given the resolved tile-based world layout from Stage 2.3, we generate context-aware visual specifications for each entity. For every entity, an LLM constructs a text prompt conditioned on the local area narrative (Stage 2.1), the entity’s semantic role, and global stylistic constraints. These prompts are passed to a text-to-3D asset generator (Hunyuan3D 2.5 Team (2025)), which produces mesh assets in .glb format independently of their final placement, enabling asset reuse across worlds. The asset generation model is treated as a black-box renderer and does not influence world structure, layout, or traversal logic.

**Orientation and Frontage Alignment.** Because text-to-3D models produce meshes with arbitrary orientation, we resolve asset orientation explicitly during world assembly. Inspired by image-grid based visual reasoning by Kim et al. (2024), we render canonical front, back, left, and right views of each generated asset and arrange them into an image grid. A vision-language model<sup>1</sup> (Dubey et al., 2024) identifies the semantic front-facing view, which is then used to orient the asset. For entities marked with `needs_frontage` (Stage 2.3), the asset’s forward direction is aligned with the normal of the nearest road edge or frontage spur tile; other entities are assigned a default or randomly sampled orientation consistent with area-level variation.

<sup>1</sup>We use meta-llama/Llama-3.2-11B-Vision-Instruct for vision-language inference

162  
163

## 2.5 STAGE 5: WORLD ASSEMBLY AND ENVIRONMENT INSTANTIATION

164  
165  
166

In the final stage, the resolved tile-based world layout and generated 3D assets are assembled into an executable environment. The global tilemap is instantiated using fixed textures for terrain and road surfaces, and each asset is placed at its assigned tile coordinates with the resolved orientation.

167  
168  
169  
170  
171

We implement this assembly step in the Godot game engine, chosen for its strong scriptability and headless execution support, which allows fully automated world construction without interactive rendering. For reproducibility, we provide the exact LLM/VLM prompt templates and the deterministic compilation details in Appendix A.1 and Appendix A.2, respectively.

172  
173

## 3 ANALYSIS &amp; DISCUSSION

174  
175  
176  
177  
178  
179

Because world structure in our system is determined algorithmically rather than learned, direct comparison to benchmark-driven predictive models would be misleading. Instead, we evaluate invariant properties of the world model using empirical metrics computed directly from the compiled representation produced at the end of Stage 2.3, where spatial layout, entity placement, and traversal structure are fully resolved.

180  
181  
182  
183  
184  
185  
186  
187  
188

Table 1 reports structural validity, functional traversability, and systems-level metrics computed from the compiled world representation. Representative qualitative examples corresponding to these statistics are shown in Appendix A.3. Structural validity is assessed via area overlap and entity collision indicators, while functional traversability is evaluated using road-network diagnostics and, critically, door reachability, which measures whether all inter-area traversal interfaces are reachable from a designated anchor. Across 20 narrative prompts, the pipeline produces collision-free, overlap-free worlds with 100% door reachability, indicating that the resulting environments are consistently executable and navigable; terminal connectivity is reported as a diagnostic and does not affect inter-area traversability when door reachability is satisfied.

189  
190  
191  
192  
193  
194  
195

Metric	Value
Collision-free worlds (%)	100.0
Area overlap-free worlds (%)	100.0
Mean terminal connectivity (diagnostic)	0.786
<b>Door reachability rate (%)</b>	<b>100.0</b>
Mean generation time (s)	$51.2 \pm 9.2$
GPU required	No (API-based inference)

196  
197  
198  
199

Table 1: Structural and functional metrics over 20 prompts. Door reachability (fraction of gates reachable from an anchor via roads) is the primary success criterion.

200  
201  
202  
203  
204  
205  
206  
207  
208

**Modeling assumptions and design scope.** The current implementation adopts several simplifying assumptions to enable deterministic compilation, including flat terrain within each area and the absence of large natural features such as rivers, lakes, coastlines, or elevation changes (small man-made water features are treated as standard entities). World structure is intentionally decomposed into modular components—including area layout, connectivity, entity placement, and asset grounding—rather than modeled within a single monolithic representation. These assumptions allow spatial layout and traversability to be resolved reproducibly within a discrete grid representation, while preserving extensibility to richer terrain variation or additional decorative assets through future modules.

209  
210  
211  
212  
213  
214  
215

**Limitations and future directions.** While the framework enables deterministic construction of coherent and traversable worlds, it does not learn environment dynamics or visual priors from data. World structure is specified through explicit rules rather than optimized to match real-world distributions or perceptual realism. An important direction for future work is to integrate this deterministic world compilation backbone with learned generative components, following hybrid approaches such as Wang et al., 2025, which combine algorithmic structure with diffusion-based image generation. Such hybrid systems could use learned models for local visual richness or stochastic variation while retaining explicit control over global structure, connectivity, and executability.

216 REFERENCES  
217

- 218 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik  
219 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling  
220 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. URL  
221 <https://arxiv.org/abs/2311.15127>.
- 222 Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes,  
223 Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative inter-  
224 active environments. *arXiv preprint arXiv:2402.15391*, 2024. URL <https://arxiv.org/abs/2402.15391>.
- 226 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, et al. The llama 3 herd of models, 2024. URL  
227 <https://arxiv.org/abs/2407.21783>.
- 229 Godot Engine Project. Godot engine, 2024. URL <https://godotengine.org>. Version 4.3.
- 231 Wonkyun Kim, Changin Choi, Wonseok Lee, and Wonjong Rhee. An image grid can be worth a  
232 video: Zero-shot video question answering using a vlm. *IEEE Access*, 12:193057–193075, 2024.  
233 doi: 10.1109/ACCESS.2024.3517625.
- 234 Zhengyang Liang, Hao He, Ceyuan Yang, and Bo Dai. Scaling laws for diffusion transformers,  
235 2024. URL <https://arxiv.org/abs/2410.08184>.
- 236 OpenAI. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- 238 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
239 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF con-  
240 ference on computer vision and pattern recognition*, pp. 10684–10695, 2022. URL <https://arxiv.org/abs/2112.10752>.
- 242 Tencent Hunyuan3D Team. Hunyuan3d 2.5: Towards high-fidelity 3d assets generation with ulti-  
243 mate details, 2025. URL <https://arxiv.org/abs/2506.16504>.
- 245 Dilin Wang et al. Worldgen: From text to traversable and interactive 3d worlds. *arXiv preprint  
arXiv:2511.16825*, 2025. URL <https://arxiv.org/abs/2511.16825>.

248 A APPENDIX  
249250 A.1 PROMPT TEMPLATES (LLM AND VLM)  
251

252 **Stage 1: Area decomposition and local narratives (LLM).** We prompt an LLM to decompose  
253 the input narrative into 3–7 outdoor areas and to produce architecture-focused descriptions intended  
254 for downstream 3D asset generation. We enforce: (i) outdoor-only areas, (ii) no time-sliced dupli-  
255 cates (e.g., day/night variants), and (iii) simplified terrain assumptions.

257 You are a creative world planner for a procedural map  
258 generator.  
259 Step 1: Plan the areas. Based on the story, identify 3–7  
260 distinct areas. – Treat areas as abstract locations, not  
261 specific moments in time. – Do NOT create separate areas  
262 that are only distinguished by time-of-day, weather, or  
263 other transient state. – Describe such changes as possible  
264 states of the SAME area.  
265 CRITICAL: ALL AREAS ARE OUTDOOR SPACES ONLY – Every area  
266 is an EXTERIOR location on the game map. – Do NOT create  
267 indoor areas (cathedral\_interior, shop\_interior, throne\_room,  
etc.). – Buildings should be described from the OUTSIDE  
only.  
268 Simplifying assumption: – Assume flat terrain; do not  
269 include large natural features such as rivers, lakes,  
cliffs, or elevation changes.

270       For each area: - Provide a unique ID (snake\_case).  
 271       - Assign a scale. - Write a narrative focusing on  
 272       ARCHITECTURAL VIBE and BUILDING STYLES: materials,  
 273       construction methods, roof/facade styles, structural  
 274       elements, and distinctive exterior details. Keep it  
 275       grounded in what players see when walking through the  
 276       outdoor area.

277 **Stage 2: Inter-area topology planning (LLM).** We prompt an LLM to produce a connected  
 278 world graph over areas, using only symbolic relative placement and discrete distance buckets. The  
 279 deterministic compiler grounds these symbolic relations into tile-grid layouts.

280       You are a world topology planner for a procedural map  
 281 generator.  
 282       Return ONLY JSON that matches the provided JSON schema  
 283 (strict). Rules: - DO NOT output coordinates, positions,  
 284 or meters. - Choose one center\_areaid. - Create  
 285 placements[] for ALL areas except the center: area\_id,  
 286 relative\_to (area\_id or "center"), dir (N/NE/E/SE/S/SW/W/NW),  
 287 dist\_bucket (near/medium/far). - Create connections[] that  
 288 forms a CONNECTED graph. - Aim for a tree (N-1 edges)  
 289 or tree + 1 loop. - Use connection kinds: trunk\_road,  
 290 road, footpath. - For each connection, specify distance  
 291 (near/medium/far). - DO NOT specify boundary edges or gate  
 292 coordinates; the compiler grounds them. Goal: produce a  
 293 connected, traversable world topology.

294 **Stage 4: Per-entity 3D asset prompt generation (LLM).** Given an entity label, tags, and (when  
 295 available) target footprint constraints from compilation, we prompt an LLM to output a short,  
 296 material-focused description suitable for text-to-3D generators.

297       You generate visual, material-focused 3D asset descriptions  
 298 for text-to-3D generators.  
 299       Hard rules: - Describe ONLY the object (no people, no  
 300 story, no actions, no camera language). - 2--4 sentences  
 301 max. - Mention: primary materials, key shapes/forms,  
 302 surface condition, and distinctive details. - If  
 303 needs\_frontage\_any is true, include a clear front with an  
 304 entrance/door orientation detail. - If placement\_dimensions  
 305 are provided, include approximate footprint and height. -  
 306 Keep it grounded in the provided context and tags.  
 307       Return ONLY valid JSON: "group": "...", "prompt": "..."

308 **Frontage selection for directional assets (VLM).** To decide which yaw view corresponds to the  
 309 main facade, we query a VLM with a  $2 \times 2$  grid of the same object rendered from four yaw directions  
 310 and request the panel containing the entrance/main facade.

311       You are given a  $2 \times 2$  grid image of the same 3D building  
 312 from four yaw directions. Define panels by position: A =  
 313 top-left, B = top-right, C = bottom-left, D = bottom-right.  
 314       Task: pick which panel shows the building's FRONT  
 315 (entrance/main facade). If the object is non-directional  
 316 or you cannot tell, answer 'NONE'.  
 317       Reply with ONLY one of: A, B, C, D, NONE.

318

319       A.2 IMPLEMENTATION / COMPIILATION DETAILS

320

321       A.2.1 STRATEGIC WORLD PLANNING DETAILS

322

323       The strategic planning stage produces a structured representation, `world_plan.json`, which encodes the semantic decomposition of the input narrative into distinct areas. Each area entry contains

324 a localized narrative description, a symbolic scale label, and a list of entities required to instantiate  
 325 the region.  
 326

327 **Symbolic Scale Categories.** Each area is assigned a discrete scale label selected from a fixed  
 328 set: tiny, small, medium, large, and huge. These symbolic scales are deterministically  
 329 mapped to predefined two-dimensional tile grid resolutions, which constrain downstream spatial  
 330 layout generation while avoiding continuous size estimation.  
 331

332 **Entity Specification.** Entities are extracted as named semantic objects without spatial coordinates.  
 333 This representation specifies *what* must exist within an area, but not *where*, ensuring that all geo-  
 334 metric decisions are deferred to later deterministic compilation stages.  
 335

### 336 A.2.2 TOPOLOGICAL GRAPH CONSTRUCTION DETAILS

337 The topological world graph (`world_graph`) encodes both the relative placement of areas and  
 338 explicit traversable connections between them. One area is designated as a global anchor, while all  
 339 other areas are positioned relative to this anchor or to previously placed areas.  
 340

341 **Symbolic Relative Placement.** Relative placement is specified symbolically using a reference  
 342 area, a compass direction, and a discrete distance category. These symbolic descriptors are deter-  
 343 ministically mapped to tile-based geometric offsets using the scale mapping defined in Appendix A,  
 344 enabling coarse spatial arrangement without metric inference from language.  
 345

346 **Inter-Area Connectivity and Gates.** Connectivity between areas is specified through named  
 347 inter-area connections. Each connection defines a traversal type (e.g., trunk road, road, footpath) and  
 348 the boundary edges of the source and destination areas on which the connection terminates. These  
 349 boundaries define abstract *gates*, which serve as symbolic traversal interfaces and are grounded to  
 350 exact tile locations during local layout compilation.  
 351

352 **Overlap Resolution.** After initial projection of symbolic placements, area bounding boxes are  
 353 tested for overlap. If conflicts are detected, inter-area distances are deterministically increased until  
 354 all layouts become disjoint. This rule-based resolution preserves global traversability while main-  
 355 taining the narrative ordering encoded in the world graph.  
 356

### 357 A.2.3 LOCAL SPATIAL LAYOUT COMPILATION DETAILS

358 Each semantic area is compiled into a local spatial layout defined over a discrete two-dimensional  
 359 tile grid. The grid resolution is determined by the symbolic scale assigned during strategic planning,  
 360 and a central anchor tile is used as the local coordinate origin.  
 361

362 **Entity Placement via Relational Constraints.** Entities are placed using relational placement con-  
 363 straints derived from semantic planning. Each constraint specifies a reference target, symbolic direc-  
 364 tion, discrete distance category, and a categorical footprint size. Entities are placed sequentially by  
 365 projecting relative offsets from the reference target; if a placement violates grid bounds or overlaps  
 366 an existing footprint, the offset distance is incrementally increased until a valid placement is found.  
 367

368 **Gate Grounding.** Traversal interfaces between areas are represented as symbolic gates and are  
 369 grounded to deterministic tile locations along the boundary edges of the area grid. Gates serve as  
 370 entry points for all intra-area road routing.  
 371

372 **Road Network Construction.** An intra-area road network is constructed directly on the tile grid to  
 373 ensure connectivity between gates and entities. Routing is performed in multiple phases, including  
 374 arterial connections to the area anchor, perimeter connections between adjacent gates, and secondary  
 375 connections to individual entities.  
 376

377 **Direction-Aware Pathfinding.** Road segments are routed using an A\*-based shortest-path search  
 over the tile grid. Each search state is augmented with an incoming direction to model turn costs and

378 directional preferences. The routing cost function incorporates penalties for sharp turns, proximity  
379 to entity footprints, and area boundaries, while encouraging reuse of existing road segments.  
380

381 **Road Materialization and Post-Processing.** Once a path is computed, all tiles along the route are  
382 marked as road tiles. Road width is enforced by lateral expansion based on road type, followed by  
383 local post-processing steps including diagonal gap filling and corner smoothing to reduce discretiza-  
384 tion artifacts.

### 386 A.3 QUALITATIVE RESULTS

388 We report qualitative outputs for 5 representative prompts (out of 20), selected to cover diverse built  
389 environments. For each prompt, we include: (i) per-area local compilation (Areas), (ii) the inter-area  
390 topological graph, and (iii) the stitched global layout with inter-area roads.

391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

432

433

434

435

436

437

438

439

440

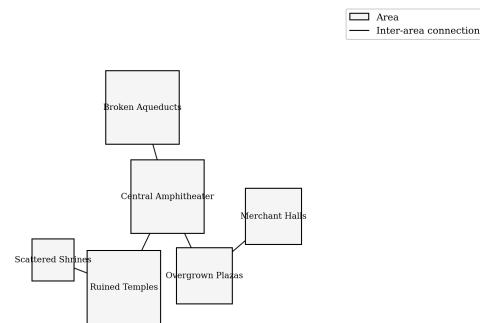
441

442

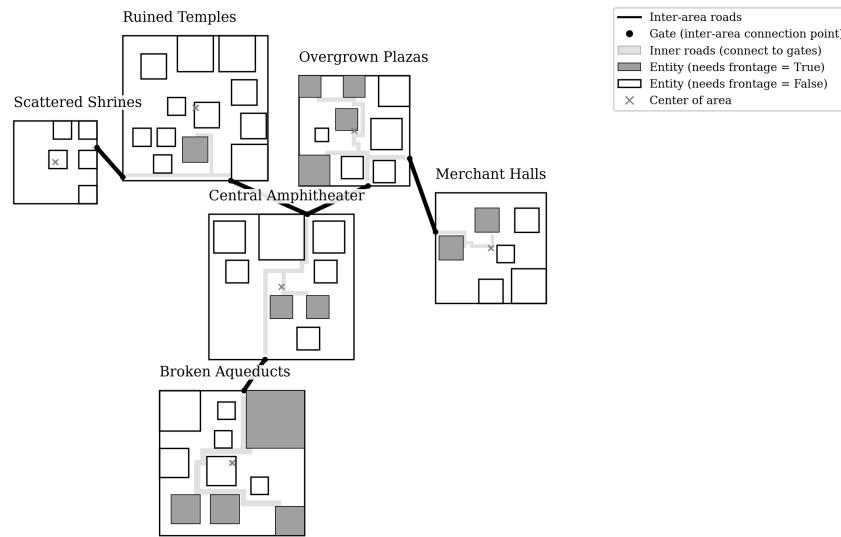
443

444

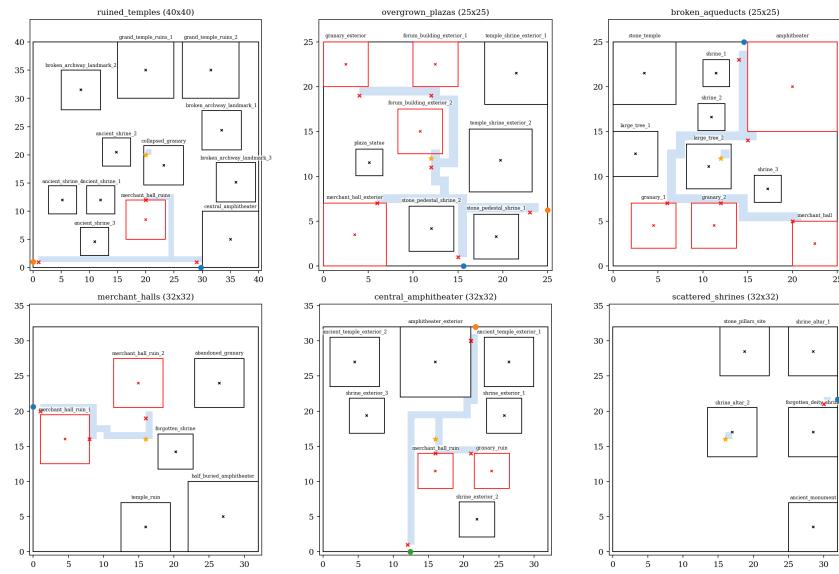
445



(a) Inter-area graph



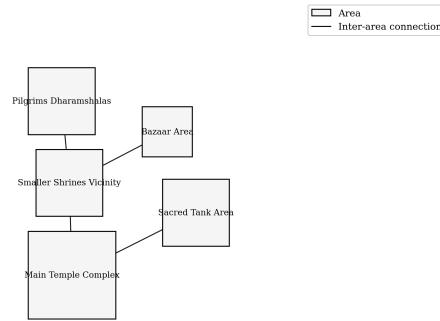
(b) Global stitched layout



(c) Per-area local compilation

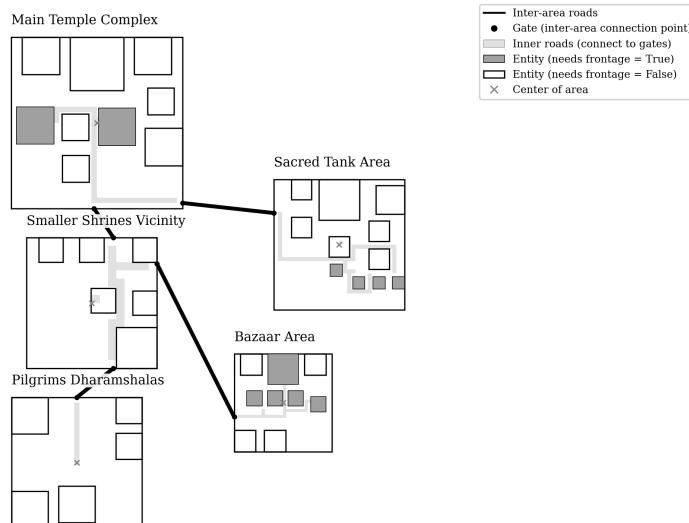
**Figure 3: Prompt 1 (world\_00).** An ancient civilization's ruins sprawl across a windswept plateau. Crumbling temples with weathered stone columns stand beside overgrown plazas. Broken aqueducts trace paths between collapsed merchant halls and abandoned granaries. A central amphitheater, half-buried in sand, hints at past gatherings. Scattered shrines mark forgotten deities, their altars still visible among the debris.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497



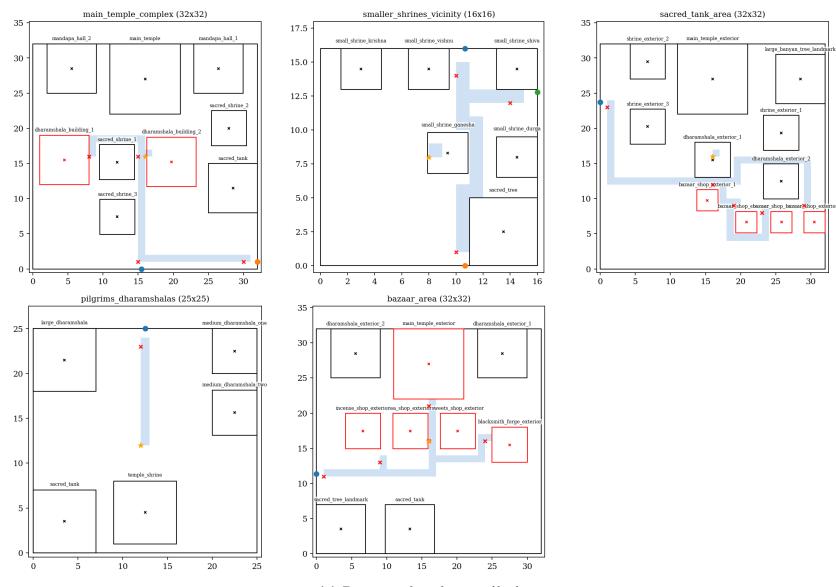
(a) Inter-area graph

498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516



(b) Global stitched layout

517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537

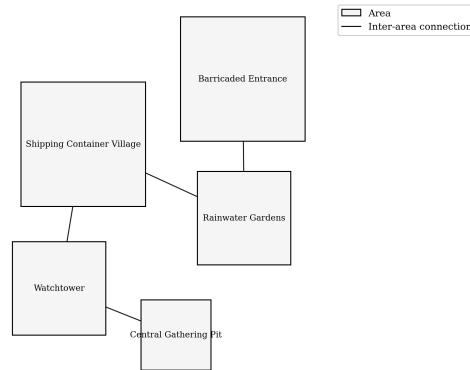


(c) Per-area local compilation

538  
539

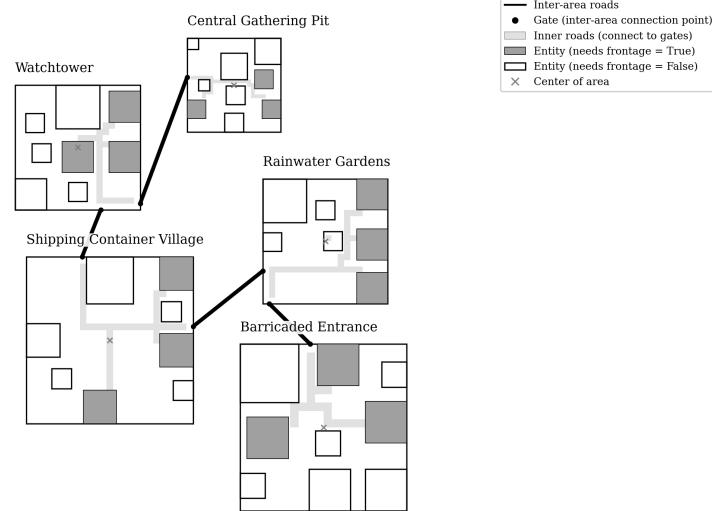
**Figure 4: Prompt (world\_03).** A sacred Indian temple complex rises from the landscape. The main temple features intricate stone carvings, gopuram towers, and mandapa halls. Smaller shrines dedicated to various deities surround the main structure. A sacred tank for ritual bathing lies to the east. Dharamshalas provide shelter for pilgrims, and a bazaar sells offerings and religious items.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552



(a) Inter-area graph

553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572



(b) Global stitched layout

573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591



(c) Per-area local compilation

592  
593

**Figure 5: Prompt (world\_07).** A survivor outpost built from the ruins of the old world. Repurposed shipping containers serve as shelters, reinforced with scrap metal. A watchtower made from scaffolding overlooks the perimeter. A central fire pit serves as the gathering place. Rainwater collectors and improvised gardens sustain the community. Barricades of wrecked vehicles protect the entrance.

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

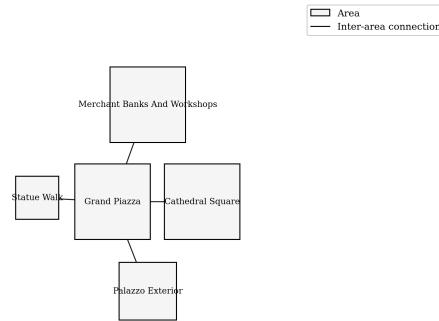
646

647



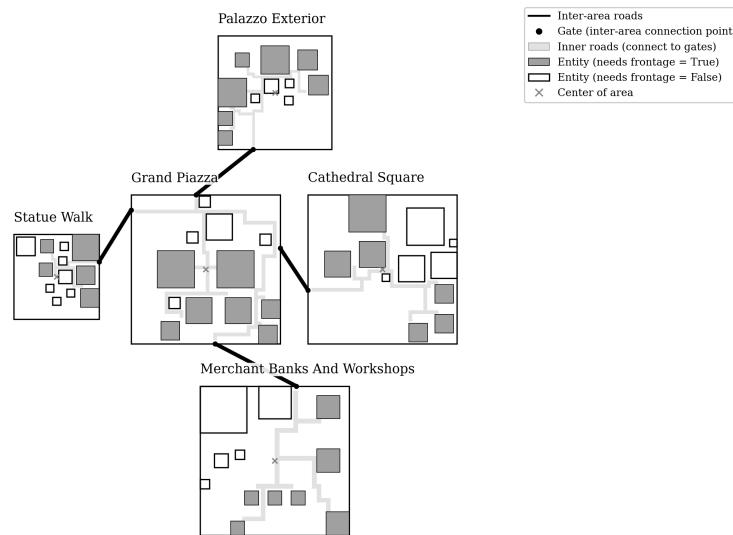
**Figure 6: Prompt (world\_10).** A Victorian industrial district shrouded in perpetual smog. Red-brick factories with tall chimneys dominate the skyline. Cobblestone streets are lined with worker tenements and public houses. A railway station connects to distant cities. Warehouses store goods from the empire. Gas lamps illuminate the fog, and a clock tower keeps time for shift changes.

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660



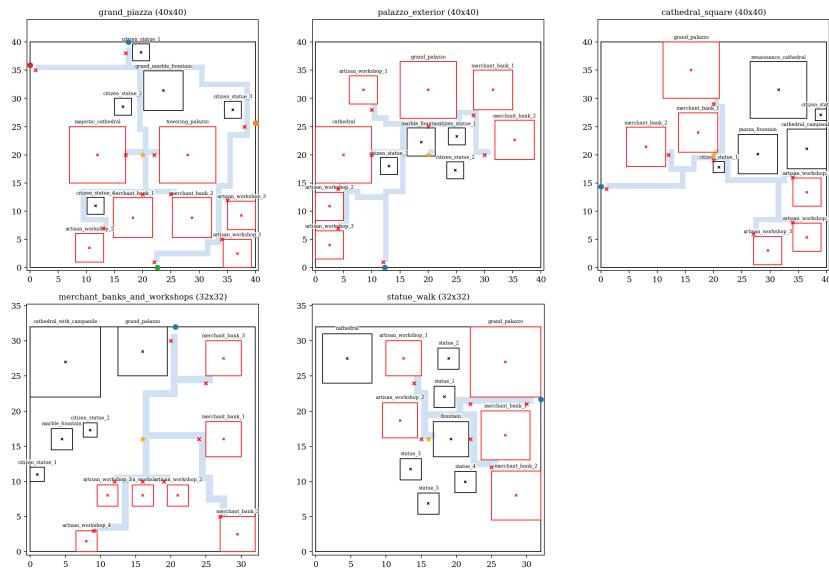
(a) Inter-area graph

661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678



(b) Global stitched layout

679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699



(c) Per-area local compilation

700  
701

**Figure 7: Prompt (world\_17).** A Renaissance Italian piazza surrounded by elegant architecture. A grand palazzo with arched colonnades faces a marble fountain. A cathedral with a bronze-domed campanile dominates one side. Merchant banks and artisan workshops line the remaining edges. Statues of famous citizens adorn the square. Cafes and tavernas welcome visitors under painted awnings.