

Symbolic Representation for Any-to-Any Generative Tasks

Anonymous CVPR submission

Paper ID 3011

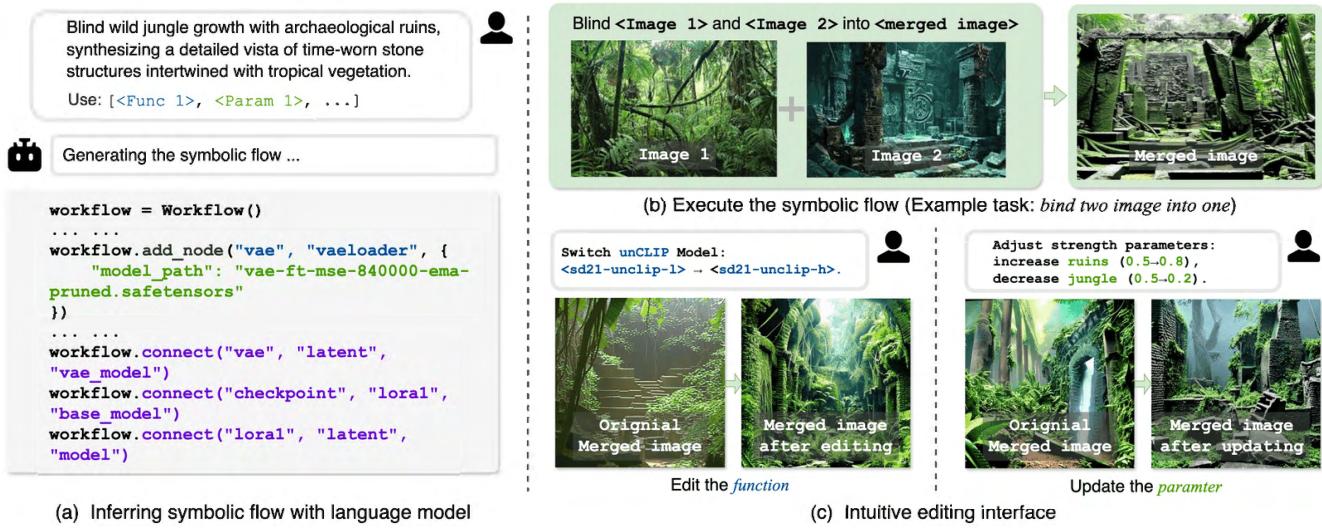


Figure 1. **A symbolic representation for Any-to-Any generative tasks.** (a) We develop a training-free inference engine that transforms natural language task descriptions into executable symbolic flow comprising *functions*, *parameters*, and the *topology*. (b) The symbolic flow allows executing generative tasks as programs. Example task is mentioned in the first sentence of Sec. 1 (c) Both *functions* and *parameters* can be easily modified to customize the generation process and the output style.

Abstract

We propose a symbolic generative task descriptive language and inference engine, capable of representing arbitrary multimodal tasks as symbolic flows. The inference engine maps natural language instructions to symbolic flow, eliminating the need for task-specific training. Conventional generative models rely heavily on large-scale training and implicit neural representation to learn cross-modal mappings, which demands extensive computational resources and restricts expandability. In this paper, we propose an explicit symbolic task descriptive language, comprising three types of primitives: *functions*, *parameters*, and *topological logic*. Using a pre-trained language model to infer symbolic workflows in a training-free manner, our framework successfully performs over 12 multimodal generative tasks based on user instructions, demonstrating enhanced efficiency and flexibility. Extensive experiments demonstrate

that our approach can generate multimodal content competitive with, and often surpassing, that of previous state-of-the-art unified models, while offering robust interruptibility and editability. We believe that symbolic task representations are capable of cost-effectively expanding the boundaries of generative AI capabilities. All code and results are available in the Supplementary Materials.

1. Introduction

“Blending the wild growth of a jungle with the mystique of ancient ruins into a brand-new scene would be stunning,” your artist friend mused. “And if we could transform the photographic image into a video, overlayed with my audio recording of birds chirping and the soft murmur of flowing water—it would create a truly dreamlike sensory experience.” This raises an interesting question:

017
018
019
020
021
022
023

024
025
026
027
028
029
030
031

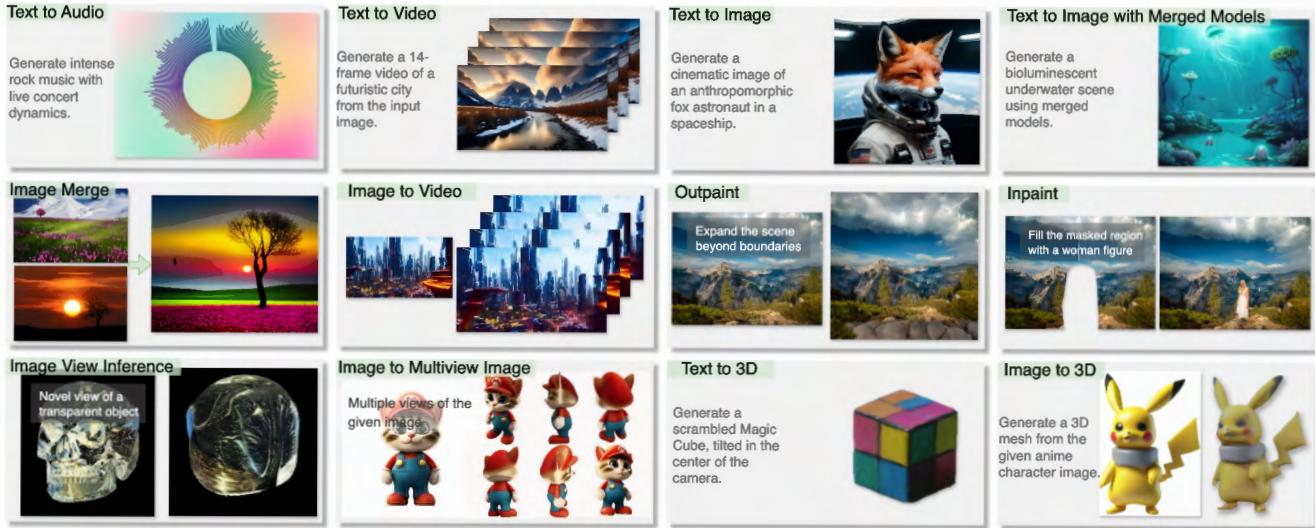


Figure 2. **The Any-to-Any generative model.** Our model demonstrates the capability to handle **any-to-any generative tasks** across various modalities, including text, images, videos, audio, and 3D content. It supports flexible transformations such as converting image to video, generating 3D models from images, or synthesizing audio from textual prompts. Formally, any-to-any generative tasks refer to generating outputs in any desired modality from inputs in any other modality, all guided by natural language instructions [42]

032 how can we design a *unified model* capable of seamlessly
 033 handling generative tasks across any combination of input
 034 and output modalities (“**any-to-any** generative tasks”, as
 035 shown in Figure 2), guided by natural language instruc-
 036 tions [12, 25, 26, 42, 49]? The workflow for executing
 037 this task comprises several essential processes [12, 39, 49].
 038 First, the system imports two images and encodes them to
 039 extract their latent features. Then, taking these features as
 040 conditioning inputs, it combines them based on the user-
 041 specified blending strength and re-synthesizes the blended
 042 latent representation onto a blank latent canvas. Finally, the
 043 system decodes this latent representation into a viewable
 044 image.

045 Current approaches for any-to-any generative tasks typ-
 046 ically fall into two paradigms: *Implicit neural modeling*
 047 and *agentic approaches*. Implicit neural modeling ap-
 048 proaches directly learn a neural representation from mass
 049 training data [25, 26, 26, 31, 40, 41, 54]. While offering
 050 simplicity in representing multimodal information, their ex-
 051 tensibility is constrained by the scope of the training data.
 052 They struggle to handle rare or unanticipated tasks—such
 053 as the image blending example in Figure 1, if such cases
 054 are not accounted for during training. Moreover, their re-
 055 liance on implicit neural representations makes them non-
 056 interruptible, leaving them ill-equipped to manage com-
 057 plex, multi-step workflows. Agentic approaches rely on
 058 sophisticated multi-agent coordination and tool orchestra-
 059 tion [12, 13, 27, 33, 38, 39], which introduces system in-
 060 stability and operational overhead in their decision-making
 061 process. While powerful, these approaches lack a unified

062 formal representation of tasks and fail to capture their in-
 063 herent compositional nature. Our experiments reveal that
 064 complex agent designs do not necessarily outperform sim-
 065 pler ones, motivating us to explore an alternative direction:
 066 focusing on *unified task representations* and *language*
 067 *model-friendly interfaces* that enable direct task specifica-
 068 tion.

069 Examining the image-blending example reveals three
 070 fundamental components essential for executing generative
 071 tasks. At its core are distinct *functions*—computational op-
 072 erations such as image encoding, conditioning, and blend-
 073 ing that transform inputs into desired outputs. Each func-
 074 tion’s behavior is shaped by *parameters*, such as the blend-
 075 ing strength and re-synthesis intensity, which fine-tune the
 076 operation to meet specific requirements. These functions
 077 do not operate in isolation; their *topology*, or interconnected
 078 relationships, form a cohesive workflow that guides the pro-
 079 gression from input to output. These three components,
 080 functions, parameters, and topology, together enable the ef-
 081 fective execution of complex generative tasks. Based on
 082 these insights, we propose *A-LANGUAGE*, a formal rep-
 083 resentation that systematically captures these three essential
 084 components of generative tasks. In *A-LANGUAGE*, *func-*
 085 *tion* specifies the core computational operations, enabling
 086 the system to precisely identify and execute required trans-
 087 formations. *parameter* provides fine-grained control over
 088 each operation’s behavior, allowing users to adapt functions
 089 to specific task requirements. *topology* formalizes the work-
 090 flow structure, defining how functions interact and com-
 091 bine to accomplish complex generative goals. Through this

092 three-component abstraction, \mathcal{A} -LANGUAGE enables flexi-
093 ble yet structured orchestration of generative tasks.

094 Alongside the symbolic generative task language, we in-
095 troduce a ***training-free inference engine*** that utilizes a pre-
096 trained language model (LM) as its foundation to derive a
097 symbolic representation from input instructions and a design-
098 ated key function. Initially, the pre-trained LM identifies a
099 comprehensive function set and parameter set from the nat-
100 ural language instruction, forming an initial functional and
101 parametric structure. With this set of functions, we then pre-
102 dict the topology, outlining the dependencies among func-
103 tions to form the complete symbolic representation. We also
104 implement a refinement module, an iterative process acti-
105 vated upon any inference failure, enabling immediate cor-
106 rections to resolve issues. Together, the \mathcal{A} -LANGUAGE, the
107 inference engine, and the refinement module led to a high-
108 quality system that provides flexible and precise workflow-
109 building capabilities.

110 Experimentally, we constructed a dataset of 120 real-
111 world generative tasks spanning 12 task categories and val-
112 idated the effectiveness of our approach through user stud-
113 ies and executability evaluations. The results demonstrate
114 that our symbolic model is competitive with or outperforms
115 state-of-the-art multimodal generative models in task gen-
116 eralization, output quality, and editing flexibility. Addi-
117 tionally, our experiments investigated the impact of syntax
118 choices on the quality of symbolic flow generated by LMs.
119 Our contributions are three-fold:

- 120 • A unified symbolic representation, the \mathcal{A} -LANGUAGE,
121 that systematically decomposes **any** generative task into
122 three core components: ***function*** for atomic operations,
123 ***parameter*** for behavioral control, and ***topology*** for sym-
124 bolic flow structure.
- 125 • A ***training-free inference engine*** that leverages pre-
126 trained LMs to automatically convert natural language
127 instructions into symbolic representations for executable
128 workflows.
- 129 • Empirical validation demonstrates that it excels in gen-
130 eralizability, modifiability, and providing an exceptional
131 user experience.

132 2. Related work

133 2.1. Unified multi-modal framework

134 Recent years have witnessed remarkable advances in large
135 language models (LLMs), which have demonstrated excep-
136 tional capabilities across various natural language tasks,
137 from basic comprehension to complex reasoning [3, 6–
138 8, 16, 21, 24, 29–31, 43, 44]. Building on this success, mul-
139 timodal large language models (MLLMs) have extended
140 these capabilities to integrate multiple forms of input and
141 output, covering data modalities such as images, audio,
142 video, and 3D structures [1, 4, 5, 10, 14, 18–20, 22, 32, 34–

37, 46, 47, 50–53, 55]. The field has progressed from iso-
143 lated single-modality models to sophisticated any-to-any
144 frameworks [25, 26, 28, 31, 40, 41, 54] that can handle
145 diverse input-output combinations within a single model
146 architecture. However, these unified multimodal frame-
147 works face significant challenges in practice. The scarcity
148 of high-quality, diverse multimodal datasets remains a funda-
149 mental bottleneck, particularly for complex cross-modal
150 tasks. Moreover, different modalities often require distinct
151 processing approaches and representations, making it chal-
152 lenging to achieve optimal performance across all possible
153 modality combinations in a single model. The need to align
154 disparate modalities into a coherent unified representation
155 while preserving their unique characteristics continues to
156 be a core challenge in advancing these frameworks.

157 2.2. Workflow synthesis

158 Workflow synthesis [2, 15, 17] seeks to generate executable
159 sequences of operations for complex tasks by coordinat-
160 ing AI models and resources, particularly in generative AI,
161 where tasks often require sophisticated combinations of in-
162 ference, parameters, and logic. Traditional methods using
163 neural modules or predefined operations struggle with the
164 open-ended nature of modern AI tasks. Recent advances
165 like HuggingGPT [39] leverage large language models for
166 task planning and model coordination, VISPROG [12] em-
167 ploys neuro-symbolic approaches for programmatic task
168 decomposition, and GenAgent [49] uses multi-agent col-
169 laboration to build workflows step by step. Despite their
170 differences, these approaches highlight the need for flexi-
171 ble, interpretable representations. Our work advances this
172 field by proposing a unified symbolic framework for de-
173 scribing and executing generative tasks, balancing expres-
174 siveness and practicality.

175 3. \mathcal{A} -Language

176 We introduce \mathcal{A} -LANGUAGE, a symbolic representation
177 that bridges the gap between natural language task descrip-
178 tions and executable workflows for any-to-any generative
179 tasks. Unlike previous unified multimodal approaches de-
180 pending on ***implicit neural representations*** and ***intensive***
181 ***training***, our \mathcal{A} -LANGUAGE provides an ***explicit symbolic***
182 ***representation***, allowing a ***training-free*** execution.

183 3.1. Formulation

184 Fundamentally, \mathcal{A} -LANGUAGE formalizes any generative
185 task t as a triple:

$$\Omega(t) := (\mathcal{F}, \Phi, \mathcal{T}). \quad 187$$

188 This unified formulation decomposes any generative task
189 into its essential constituents: the computational ***functions***
190 \mathcal{F} , their corresponding ***parameters*** Φ , and the ***topological***

191 **structure** \mathcal{T} that elucidates their interrelations and data flow
 192 dynamics.

193 **Function** The function set is defined as $\mathcal{F} =$
 194 $\{f_1, f_2, \dots, f_n\}$, where $n \in \mathbb{N}$, which represents atomic
 195 computational units. Each function takes both input data
 196 and parameters to produce outputs, formally defined as:

$$197 \quad f_i : \mathcal{I}_i \times \phi_i \rightarrow \mathcal{O}_i,$$

198 where \mathcal{I}_i defines its input space, ϕ_i represents its parameter
 199 configuration, and \mathcal{O}_i specifies its output space. The input
 200 and output spaces \mathcal{I}_i and \mathcal{O}_i represent either simple scalar
 201 values or composite data structures of arbitrary modalities,
 202 allowing functions to process multiple inputs and generate
 203 multiple outputs. For example, an image blending function
 204 might accept two image inputs and produce both a
 205 blended result and an attention mask. When functions are
 206 connected, their inputs and outputs can be partially mapped,
 207 providing flexibility in constructing complex paths.

208 **Parameter** The parameter space $\Phi = \{\phi_{f_1}, \phi_{f_2}, \dots, \phi_{f_n}\}$
 209 encompasses configurations that modify function behaviors,
 210 where each ϕ_{f_i} represents the parameter space for
 211 function f_i . Parameters must be fully specified before func-
 212 tion execution to ensure deterministic behavior. The param-
 213 eter space is independent of the input space, enabling func-
 214 tions to exhibit different behaviors while processing identi-
 215 cal inputs.

216 **Topology** The topology set $\mathcal{T} = \{d_1, d_2, \dots, d_m\}$ defines
 217 the precise data flows between functions, where each d_k at
 218 the finest granularity specifies a single directed connection
 219 from a specific output of one function to a specific input
 220 of another function. Specifically, d_k is defined as a tuple
 221 representing an individual data flow from the output of a
 222 source function to the input of a target function. Formally:

$$223 \quad d_k = (f_j, y_j) \rightarrow (f_i, x_i) \mid y_j \in \mathcal{O}_j, x_i \in \mathcal{I}_i$$

224 where f_j and f_i denote the source and target functions, re-
 225 spectively. y_j refers to a specific output produced by func-
 226 tion f_j , while x_i corresponds to a specific input required by
 227 function f_i . Thus, each d_k encapsulates the transfer of data
 228 from a designated output of one function to a designated
 229 input of another, allowing for precise tracking of data flow
 230 through the system.

231 **Symbolic flow** The symbolic flow emerges from the in-
 232 teraction of **functions**, **parameters**, and **topological logic**,
 233 formalizing the complete generative process:

$$234 \quad \mathcal{S} = \{(f_i, \phi_{f_i}, D_i) \mid f_i \in \mathcal{F}\},$$

where D_i is the set of all data flows d_k in \mathcal{T} that target
 235 function f_i :

$$236 \quad D_i = \{(f_j, y_j) \rightarrow (f_i, x_i) \mid f_j \in \mathcal{F}, y_j \in \mathcal{O}_j, x_i \in \mathcal{I}_i\}.$$

237 Each element in the symbolic flow specifies a function, its
 238 parameter configuration, and its incoming directed connec-
 239 tions. Specifically, for each function f_i , D_i contains tu-
 240 ples that map specific outputs of predecessor functions to
 241 specific inputs of f_i . This fine-grained formulation cap-
 242 tures how computation progresses through the system, with
 243 functions receiving their required inputs from designated
 244 outputs of antecedent functions and parameter configura-
 245 tions from the parameter space. Through this unified and
 246 detailed representation, *A-LANGUAGE* can express diverse
 247 and complex generative tasks. m

3.2. Syntax styles

249 The symbolic representation $\Omega(t)$ can be expressed through
 250 multiple syntactic styles, as shown in Figure 3, each of-
 251 fering different trade-offs in expressiveness and clarity. To
 252 identify the most effective representation for large language
 253 model inference, we explore three distinct syntactic formu-
 254 lations: **declarative**, **dataflow**, and **pseudo-natural** syntax,
 255 as illustrated through concise examples in Figure 3.

257 **Declarative Syntax** Declarative Syntax [45] focuses on
 258 explicitly specifying computational components and their
 259 relationships. Functions are separately declared with pa-
 260 rameters, while connections are specified through explicit
 261 statements. This style is effective for complex workflows
 262 with reusable components, as it clearly separates compo-
 263 nent definitions (\mathcal{F}) from relationships (\mathcal{T}).

264 **Dataflow syntax** Dataflow syntax [49] emphasizes the
 265 flow of data through function compositions, where out-
 266 puts directly feed into subsequent functions. It captures
 267 topological relationships (\mathcal{T}) through the order of func-
 268 tion calls while maintaining explicit parameter specifi-
 269 cations (Φ). This style is particularly suited for linear, sequen-
 270 tial workflows.

271 **Pseudo-natural syntax** Pseudo-natural syntax [9] aims
 272 to bridge formal representations with more intuitive,
 273 language-like structures, making task specifications more
 274 accessible while maintaining mathematical rigor. This style
 275 explores a balance between precision and readability.

276 Each style retains the full expressiveness of $\Omega(t)$, but of-
 277 fers different advantages in terms of clarity and usability.
 278 The subsequent empirical analysis will evaluate which syn-
 279 tax best supports natural language inference while preserv-
 280 ing necessary formal properties.

Notation	Implementation and definition
System Components	
\mathcal{X}	List [Any] // Input data of any modality
s	str // Task description
\mathcal{C}	Dict // System constraints
$\Omega(t)$	Workflow // Complete workflow representation
Workflow Structure	
$f_i \in \mathcal{F}$	Node // Computational function
$f_i : \mathcal{I}_i \times \phi_i \rightarrow \mathcal{O}_i$	Node.forward // Function mapping with parameters
$\phi_{f_i} \in \Phi$	Dict[str, Any] // Function parameters
$d_k \in \mathcal{T}$	(Node, Any) -> (Node, Any) // Source output to target input mapping $((f_j, y_j) \rightarrow (f_i, x_i))$
Workflow Operations	
Initialize	Workflow() // Create empty workflow $\Omega(t) = (\mathcal{F}, \Phi, \mathcal{T})$
Add Node	add_node(name, type, params) // Add function f_i with parameters ϕ_{f_i}
Connect	connect(src_node, src_output, dst_node, dst_input) // Create topology $d_k : (f_j, y_j) \rightarrow (f_i, x_i)$

Table 1. System components and operations summary. A comprehensive overview of \mathcal{A} -LANGUAGE’s system components and their implementations. The upper two sections define the mathematical notations and their corresponding implementations, where the system processes input data \mathcal{X} according to task description s under constraints \mathcal{C} . Functions f_i transform inputs \mathcal{I}_i with parameters ϕ_i to outputs \mathcal{O}_i , and are connected through directed mappings d_k . The lower section demonstrates the Declarative Syntax as one example of workflow construction, showing how basic operations map to the mathematical formulation $\Omega(t) = (\mathcal{F}, \Phi, \mathcal{T})$.

```

workflow = Workflow()
.....
workflow.add_node("vae",
  "vaeloader",
  "model_path": "vae-ft-
  mse-840000-ema-
  pruned.safetensors"
)
.....
workflow.connect("vae",
  "latent", "vae_model")

```

(a) Declarative Syntax

```

vae = vaeloader(
  model_path="vae-ft-
  mse-840000-ema-
  pruned.safetensors"
)
...
vae = vae_model(
  latent
)

```

(b) Dataflow Syntax

```

vae is vaeloader with
the parameter of
(model path is "vae-
ft-mse-840000-ema-
pruned.safetensors")
...
vae is vae model with
the parameter of
(latent)

```

(c) Pseudo-natural Syntax

Figure 3. Syntax comparison. We implement our symbolic representation using three different styles of domain-specific languages (DSLs). (a) The declarative syntax registers all components into the workflow. (b) The dataflow syntax emphasizes the direction of data flow. (c) The pseudo-natural syntax mimics human language expression.

4. Inferring via pre-trained language model

The diversity and complexity of generative tasks necessitate a flexible and robust approach to transforming high-level task specifications into executable symbolic flows. As illustrated in Figure 4, we propose utilizing LMs as inference engines to generate task-specific symbolic representations, with Figure 5 demonstrating the complete pipeline from natural language description to executable workflow. This enables any-to-any transformations across different modalities and task types.

Given a set of inputs \mathcal{X} of arbitrary modalities, a task description s , and a set of constraints \mathcal{C} , our inference framework generates a complete symbolic representation $\Omega(t)$. As illustrated in Figure 4, our framework leverages a pre-trained language model to infer both the computational components and their topology from natural language descriptions. This process can be formalized as:

$$\mathcal{M} : (\mathcal{X}, s, \mathcal{C}) \rightarrow \Omega(t),$$

where \mathcal{X} represents any combination of inputs such as images, text, audio, or other modalities, s describes the desired transformation, and \mathcal{C} represents a set of constraints, which typically specifying information such as available functions, specific parameter choices, valid parameter ranges, and model compatibility. These constraints are essential for ensuring that the generated symbolic flow is not only theoretically sound but also practically executable within the given computational environment. Specifically, we divide the inference into three main steps:

Component inference The first stage of our framework focuses on determining the necessary computational components. Given the input specifications and constraints, the LM identifies the required functions and their parameters:

$$\psi_1 : (\mathcal{X}, s, \mathcal{C}) \rightarrow (\mathcal{F}, \Phi).$$

This process accounts for both the explicit requirements of the task and any implicit dependencies, ensuring that selected functions are available within \mathcal{C} .

Topology construction The second stage focuses on establishing relationships between the identified components to form a coherent computational flow:

$$\psi_2 : (\mathcal{X}, s, \mathcal{C}, \mathcal{F}, \Phi) \rightarrow \mathcal{T}.$$

In this phase, the LM evaluates how the outputs of one function can serve as inputs to another, ensuring that these connections are executable and comply with the constraints defined in \mathcal{C} . This construction guarantees that data flows seamlessly through the system in a manner consistent with our unified formulation.

299
300
301
302
303
304
305
306
307
308

309
310
311
312

317
318
319

321
322
323
324
325
326

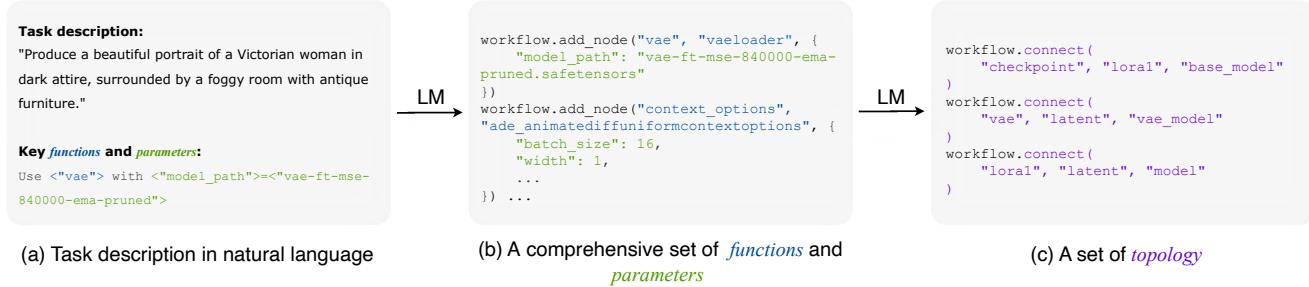


Figure 4. **Inferring symbolic flow with pre-trained language model (LM).** Beginning with (a) a natural language task description and key functions and parameters, we leverage LM to infer (b) a comprehensive set of functions and parameters. We then integrate (a) and (b) to deduce the (c) topology. If compilation or execution fails, all information is aggregated for further refinement (Sec. 4).

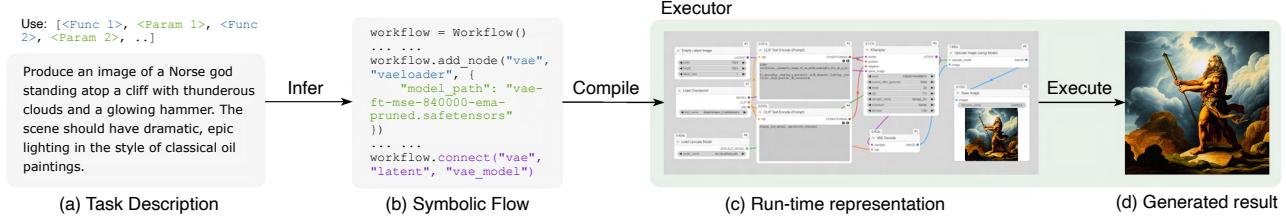


Figure 5. **Demonstration of the inference and execution.** The inference framework translates a natural language task description into an executable symbolic representation. This symbolic representation is then compiled and executed through a workflow executor to perform the desired transformation. See appendix for details.

Iterative refinement The generated symbolic flow undergoes an iterative refinement process to ensure correctness and executability. We define this refinement as:

$$\Omega_{i+1}(t) = R(\Omega_i(t), \epsilon_i),$$

where R represents the refinement operator and ϵ_i captures any detected issues in iteration i . To prevent endless loops, a maximum number of iterations can be set. During each iteration, the LM analyzes error signals and adjusts the symbolic flow accordingly, either by modifying function parameters, adding missing components, or restructuring topological connections. This iterative process continues until a valid symbolic flow is achieved that satisfies all constraints in \mathcal{C} or the maximum iteration count is reached.

The combination of LM-based inference and iterative refinement enables our framework to handle diverse transformation tasks while maintaining robustness and generality. By leveraging the LM’s reasoning capabilities and incorporating explicit constraints, we bridge the gap between high-level task descriptions and executable symbolic flows, providing a flexible foundation for any-to-any transformations.

5. Experiments

5.1. Setup

Prompt suite We collected a diverse set of 120 generative tasks from real-world applications to comprehensively evaluate our approach (see Appendix for the complete task

list). These tasks are categorized into 12 general groups, each comprising 10 distinct instances. See Appendix for details.

Table 2. **Comparison of the average rankings** between outcome quality and task-outcome alignment rankings (\downarrow). We primarily compared ***neural representing, training-dependent modeling*** [11, 23, 26, 48] and our ***symbolic representing, training-free modeling***. Each method was ranked on a scale starting from 1, with 1 denoting the best-performing approach. “U-IO 2” denotes “Unified-IO”, “I-2-3D” denotes “Image to 3D Mesh”, “T2M” denotes “Text to Mesh”.

Method	Inpaint	Outpaint	Img merge	NVS Merge	model I-2-3D	
Show-o [48]	1.6	1.4	X	X	X	X
SEED-X [11]	X	X	1.2	X	X	X
LWM [23]	X	X	X	X	X	X
U-IO 2 [26]	-	X	-	X	X	X
Ours	1.4	1.6	1.8	1.0	1.0	1.0
Method	T2I	T2A	Multi-view img	I2V	T2M	T2V
Show-o [48]	2.8	X	X	X	X	X
SEED-X [11]	2.0	X	X	X	X	X
LWM [23]	4.2	X	X	X	X	X
U-IO 2 [26]	4.5	2.0	-	-	X	X
Ours	1.5	1.0	1.0	1.0	1.0	1.0

Metric ① For execution evaluation, we first evaluated the single-run ***pass rate (Pass@1)*** of compilation and execution, following Xue *et al.* [49]. ② For outcome quality and instruction-following, we conducted a systematic user study

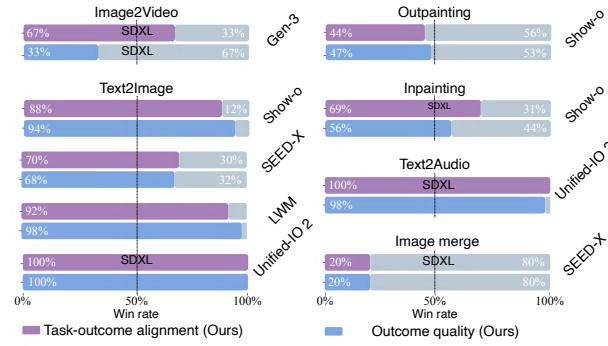


Figure 6. **Comparison of our win rates** with the state-of-the-art unified multimodal models.

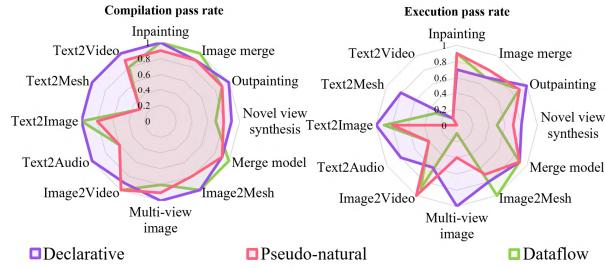


Figure 7. **Comparison of syntax styles.** Metric: Pass@1 (\uparrow). See Appendix for details.

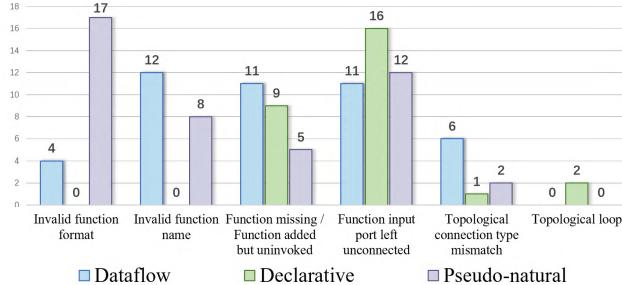


Figure 8. **Comparative error distribution** for dataflow, declarative, and pseudo-natural syntax styles, illustrating six types of errors occur when testing on the 120 generative tasks.

with five annotators who ranked outputs from all frameworks for comparison, the metrics are following:

- **Text-outcome alignment:** We measured the degree of correspondence between generated outputs and their intended task specifications. Higher alignment scores indicated closer matches between system outputs and expected results based on input requirements.
- **Outcome quality:** We assessed generated outputs based on three criteria: aesthetic appeal, structural coherence, and technical quality. This metric encompassed visual clarity, presentation effectiveness, and adherence to task-specific quality standards.

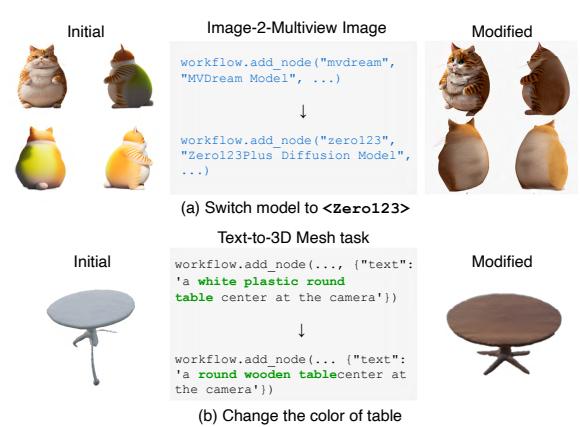


Figure 9. **Symbolic Flow Editing.** We present examples of modifying (a) *functions*, where users can directly change models by editing code to achieve desired effects, and (b) *parameters*, such as adjusting textual prompts (treated as a type of parameter) to alter the color of 3D assets.

- **Average rank:** We computed this metric by first ranking each model’s performance on text-outcome alignment and outcome quality for individual samples, then calculating the mean rank across all tasks.
- **Win rate:** A “win” is recorded when our method ranks higher than a competitor for a given sample. The win rate represented the percentage of successful comparisons, serving as a measure of relative performance advantage.

Table 3. **Agentic design [49] vs. symbolic inference (Ours).** We calculate the average pass rate (Pass@1, \uparrow) on compilation and execution. Results are averaged across 120 generative tasks.

Method	Compilation	Execution
GenAgent [49]	0.84	0.63
Ours	0.97	0.77

Baselines ① **Agentic framework:** We selected GenAgent [49] as our primary baseline method. To ensure fairness, we augmented GenAgent [49] with key functions and parameters as additional input, and increased the maximum refinement iterations to 3. ② **Unified multimodal models:** We also compared against the state-of-the-art unified multimodal approaches. In the Text to Image and Inpaint tasks, the Show-o model [48] has a guidance scale of 1.75 and 16 time steps. For Outpaint tasks, we set both left and right expansion degrees to 1. The SEED-x model [11] was configured with a maximum output token count of 1024 and a maximum of 3 history rounds. We enabled three specific options: forced image generation, forced bounding boxes, and forced image optimization. ③ **Commercial generative model:** The Gen-3 video generation model [37] was configured with 720p resolution (1280×768 aspect ratio), using random seed and a video length of 5 seconds.

371
372
373
374
375
376
377
378

379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395

396 **Implementation details** Following Gupta *et al.* [12], we
397 implemented in-context learning to prompt the LM with
398 syntax and logical guidance. Specifically, we performed
399 Retrieval-Augmented Generation (RAG) based on the task
400 description, retrieving three most relevant programs as ref-
401 erences. We curated a reference program database con-
402 taining 16 distinct programs, ensuring no overlap with
403 the target evaluation tasks. All experiments were con-
404 ducted on a single L4 GPU (24GB), with 1TB external
405 storage, running on a Debian 11 server. ComfyUI
406 served as the back-end for code execution. We used GPT-
407 4o (gpt-4o-2024-08-06) as the inference engine and
408 text-embedding-3-large as the embedding model.

409 5.2. Main results

410 **Comparative performance in user study** Our sym-
411 bolic model consistently outperforms state-of-the-art uni-
412 fied models in both text-outcome alignment and result qual-
413 ity across multiple generative tasks. In the user study
414 involving five experienced participants, our model was
415 evaluated against Show-o [48], SEED-X [11], LVM [23],
416 Unified-IO [26], and the commercial Gen-3 [37]. As illus-
417 trated in Figure 6, our approach achieved a 94% win rate
418 against Show-o [48] and 98% against LVM [23] in Text
419 to Image tasks. Notably, in Image2Video generation, our
420 model surpassed the commercial Gen-3 with a 67% win rate
421 in text-outcome alignment. Additionally, for Text to Audio,
422 our model attained a 100% win rate in alignment and 98%
423 in quality against Unified-IO [26], underscoring its superior
424 performance across diverse applications. See Appendix for
425 the visualization results.

426 **Is complex agentic design necessary?** As shown in Ta-
427 ble 3, simpler, symbolic approaches can achieve higher suc-
428 cess rates for straightforward tasks without the complexities
429 and costs associated with agentic designs. Unlike GenA-
430 gent [49], which employs multi-step planning and actions
431 that can amplify errors and increase computational costs,
432 our symbolic method maintains simplicity and clarity. This
433 reduction in complexity leads to higher success rates in sim-
434 ple tasks by minimizing error propagation and lowering ex-
435 ecution costs. However, for more intricate workflows, in-
436 tegrating symbolic representations with agentic strategies
437 may offer enhanced flexibility and performance, suggesting
438 a potential hybrid approach for future research. See Ap-
439 pendix for details.

440 **Representation: neural or symbolic?** Our symbolic
441 model outperforms neural models in task generality and
442 output quality without additional training. Table 2 high-
443 lights that our symbolic approach successfully handles all
444 120 generative tasks, including complex categories such as

3D and video generation. In contrast, neural models are lim-
445 ited by their reliance on extensive training data, restricting
446 their ability to manage diverse and complex tasks. Specifi-
447 cally, our model achieves superior average ranks in most 2D
448 tasks like Inpaint, Text to Audio, and Text to Image gener-
449 ation, demonstrating its enhanced adaptability and perfor-
450 mance over unified neural frameworks.

452 **Explicit symbolic flow editing** Our symbolic represen-
453 tation enables precise and effective control over distinct
454 stages of generative tasks, thus paving the way for the re-
455 alization of more complex tasks. Figure 9 illustrates exam-
456 ples of modifying *function* (model) and *parameter* (textual
457 prompt), respectively. By applying explicit program modifi-
458 cations, control over the image generation process is given.
459 See Appendix for more examples.

460 **Error analysis: What constitutes an LM-friendly syn-
461 tax style?** A balance between human readability and for-
462 mat correctness is essential for enhancing language model
463 performance, with structural rigidity impacting topological
464 clarity. Upon analysis of the reasoning processes of the 120
465 test tasks in Figure 8, we identified two main takeaways.

- **❶ Human readability vs. format correctness:** Higher
466 readability in language design correlates with increased
467 format errors. Pseudo-natural language formats exhibited
468 17 instances of invalid code formats, compared to 4 in
469 dataflow and none in declarative styles. This indicates
470 that while readability facilitates human understanding, it
471 can hinder precise format adherence by language models.
- **❷ Structural rigidity vs. topological clarity:** Struc-
472 turally rigid and highly modular languages, such as our
473 declarative syntax, tend to introduce topological gaps and
474 connection errors, with 9 instances of missing or unin-
475 voked functions and 16 unlinked input ports. This sug-
476 gests that increased structural complexity can challenge
477 language models in maintaining clear and accurate depen-
478 dencies between functions and ports.

481 6. Conclusion

482 We have proposed a symbolic generative task description
483 language, combined with an inference engine, provid-
484 ing a novel and efficient way to represent and execute
485 multimodal tasks without the need for task-specific
486 training. By leveraging a pre-trained large language
487 model to infer symbolic task descriptions, our approach
488 has successfully synthesized diverse multimodal tasks,
489 demonstrating its flexibility and potential to unify dif-
490 ferent generative AI capabilities. Our experiments on
491 120 tasks have shown that our framework has achieved
492 performance comparable to unified multimodal mod-
493 els, highlighting its expandability and cost-effectiveness.

495 **References**

- [1] Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing Gong. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text. *Advances in Neural Information Processing Systems*, 34:24206–24221, 2021. 3
- [2] Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural module networks, 2017. 3
- [3] Anthropic. The claudie 3 model family: Opus, sonnet, haiku, 2024. Corpus ID: 268232499. 3
- [4] James Betker, Gabriel Goh, Li Jing, † TimBrooks, Jianfeng Wang, Linjie Li, † LongOuyang, † JuntangZhuang, † JoyceLee, † YufeiGuo, † WesamManassra, † PrafullaDhariwal, † CaseyChu, † YunxinJiao, and Aditya Ramesh. Improving image generation with better captions. 3
- [5] Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. Audioltm: a language modeling approach to audio generation, 2023. 3
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 2020. 3
- [7] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. 2023.
- [8] DeepSeek-AI, Xiao Bi, Deli Chen, Guanting Chen, Shanhua Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, Huazuo Gao, Kaige Gao, Wenjun Gao, Ruiqi Ge, Kang Guan, Daya Guo, Jianzhong Guo, Guangbo Hao, Zhewen Hao, Ying He, Wenjie Hu, Panpan Huang, Erhang Li, Guowei Li, Jiashi Li, Yao Li, Y. K. Li, Wenfeng Liang, Fangyun Lin, A. X. Liu, Bo Liu, Wen Liu, Xiaodong Liu, Xin Liu, Yiyuan Liu, Haoyu Lu, Shanghai Lu, Fuli Luo, Shirong Ma, Xiaotao Nie, Tian Pei, Yishi Piao, Junjie Qiu, Hui Qu, Tongzheng Ren, Ze-hui Ren, Chong Ruan, Zhangli Sha, Zhihong Shao, Junxiao Song, Xuecheng Su, Jingxiang Sun, Yaofeng Sun, Minghui Tang, Bingxuan Wang, Peiyi Wang, Shiyu Wang, Yaohui Wang, Yongji Wang, Tong Wu, Y. Wu, Xin Xie, Zhenda Xie, Ziwei Xie, Yiliang Xiong, Hanwei Xu, R. X. Xu, Yan-hong Xu, Dejian Yang, Yuxiang You, Shuiping Yu, Xingkai Yu, B. Zhang, Haowei Zhang, Lecong Zhang, Liyue Zhang, Mingchuan Zhang, Minghua Zhang, Wentao Zhang, Yichao Zhang, Chenggang Zhao, Yao Zhao, Shangyan Zhou, Shun-feng Zhou, Qihao Zhu, and Yuheng Zou. Deepseek llm: Scaling open-source language models with longtermism, 2024. 3
- [9] Michael D Ernst. Natural language is a programming language: Applying natural language processing to software development. In *2nd Summit on Advances in Programming Languages (SNAPL 2017)*. Schloss-Dagstuhl-Leibniz-Zentrum für Informatik, 2017. 4
- [10] Han Fang, Pengfei Xiong, Luhui Xu, and Yu Chen. Clip2video: Mastering video-text retrieval via image clip. *arXiv preprint arXiv:2106.11097*, 2021. 3
- [11] Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Multimodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024. 6, 7, 8
- [12] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *CVPR*, pages 14953–14962, 2023. 2, 3, 8
- [13] Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings, 2024. 2
- [14] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers, 2022. 3
- [15] Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. Learning to reason: End-to-end module networks for visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 804–813, 2017. 3
- [16] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024. 3
- [17] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Judy Hoffman, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Inferring and executing programs for visual reasoning. In *Proceedings of the IEEE international conference on computer vision*, pages 2989–2998, 2017. 3
- [18] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation, 2023. 3
- [19] Junnan Li, Ramprasaath R Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021.
- [20] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 3
- [21] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muenninghoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, 9

- 609 Rudra Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry 666
610 Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour 667
611 Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, 668
612 Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zh- 669
613 danov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer 670
614 Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri 671
615 Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Car- 672
616 olyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contrac- 673
617 tor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine 674
618 Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas 675
619 Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 676
620 Starcoder: may the source be with you!, 2023. 3
621 [22] Xiuju Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei 677
622 Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu 678
623 Wei, et al. Oscar: Object-semantics aligned pre-training 679
624 for vision-language tasks. In *Computer Vision–ECCV 2020: 680*
625 16th European Conference, Glasgow, UK, August 23–28, 681
626 2020, Proceedings, Part XXXI 6, pages 121–137. Springer, 682
627 2020. 3
628 [23] Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. 683
629 World model on million-length video and language with 684
630 ringattention. *arXiv preprint arXiv:2307.16789*, 2024. 6, 8
631 [24] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico 685
632 Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtiar, Jiawei Liu, Yuxiang Wei, Tianyang 686
633 Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes 687
634 Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraein 688
635 Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia 689
636 Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, 690
637 Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman 691
638 Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, 692
639 Yekun Chai, Niklas Muenninghoff, Xiangru Tang, Muh- 693
640 tasham Oblokulov, Christopher Akiki, Marc Marone, Cheng- 694
641 hao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, 695
642 Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, 696
643 Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, 697
644 Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapa- 698
645 dos, Mostafa Patwary, Nima Tajbakhsh, Yacine Jernite, 699
646 Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, 700
647 Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm 701
648 de Vries. Starcoder 2 and the stack v2: The next generation, 702
649 2024. 3
650 [25] Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh 703
651 Mottaghi, and Aniruddha Kembhavi. Unified-io: A unified 704
652 model for vision, language, and multi-modal tasks, 2022. 2, 705
653 3
654 [26] Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, 706
655 Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha 707
656 Kembhavi. Unified-io 2: Scaling autoregressive multimodal 708
657 models with vision, language, audio, and action, 2023. 2, 3, 709
658 6, 8
659 [27] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei 710
660 Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 711
661 Chameleon: Plug-and-play compositional reasoning with 712
662 large language models. *Advances in Neural Information Pro- 713
663 cessing Systems*, 36, 2024. 2
664 [28] David Mizrahi, Roman Bachmann, Oğuzhan Fatih Kar, 714
665 Teresa Yeo, Mingfei Gao, Afshin Dehghan, and Amir Zamir. 715
666 4m: Massively multimodal masked modeling, 2023. 3
667 [29] OpenAI. Chatgpt: Optimizing language models for 716
668 dialogue. <http://web.archive.org/web/20230109000707/https://openai.com/blog/chatgpt/>, 2022. 3
669 [30] OpenAI et al. Gpt-4 technical report, 2024. 672
670 [31] OpenAI et al. Gpt-4o system card, 2024. 2, 3
671 [32] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion, 2022. 3
672 [33] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, 673
673 Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, 674
674 et al. Toollm: Facilitating large language models to master 675
675 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 676
676 2023. 2
677 [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 678
678 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, 679
679 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning 680
680 transferable visual models from natural language supervi- 681
681 sion. In *International conference on machine learning*, pages 682
682 8748–8763. PMLR, 2021. 3
683 [35] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, 684
684 Christine McLeavey, and Ilya Sutskever. Robust speech 685
685 recognition via large-scale weak supervision, 2022. 686
686 [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 687
687 Patrick Esser, and Björn Ommer. High-resolution image syn- 688
688 thesis with latent diffusion models, 2022. 689
689 [37] Runway. Gen-3. <https://runwayml.com/blog/introducing-gen-3-alpha/>, 2024. 3, 7, 8
690 [38] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta 691
691 Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Can- 692
692 cedda, and Thomas Scialom. Toolformer: Language models 693
693 can teach themselves to use tools, 2023. 2
694 [39] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, 695
695 Weiming Lu, and Yuetong Zhuang. Hugginggpt: Solving ai 696
696 tasks with chatgpt and its friends in hugging face, 2023. 2, 3
697 [40] Zineng Tang, Ziyi Yang, Mahmoud Khademi, Yang Liu, 698
698 Chenguang Zhu, and Mohit Bansal. Codi-2: In-context, 699
699 interleaved, and interactive any-to-any generation, 2023. 2, 3
700 [41] Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and 701
701 Mohit Bansal. Any-to-any generation via composable diffu- 702
702 sion, 2023. 2, 3
703 [42] Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and 704
704 Mohit Bansal. Any-to-any generation via composable diffu- 705
705 sion. *Advances in Neural Information Processing Systems*, 706
706 36, 2024. 2
707 [43] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier 708
708 Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste 709
709 Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aure- 710
710 lien Rodriguez, Armand Joulin, Edouard Grave, and Guilla- 711
711 laume Lample. Llama: Open and efficient foundation lan- 712
712 guage models, 2023. 3
713 [44] Hugo Touvron, Louis Martin, Kevin Stone, Peter Al- 714
714 bert, Amjad Almahairi, Yasmine Babaei, Nikolay Bash- 715
715 lykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, 716
716 Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, 717
717

- 722 Moya Chen, Guillem Cucurull, David Esiobu, Jude Fer-
723 nandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia
724 Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
725 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Vik-
726 tor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Ko-
727 renov, Punit Singh Koura, Marie-Anne Lachaux, Thibaut
728 Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning
729 Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
730 Igor Molybog, Yixin Nie, Andrew Poultion, Jeremy Reizen-
731 stein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan
732 Silva, Eric Michael Smith, Ranjan Subramanian, Xiao-
733 qing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams,
734 Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov,
735 Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan
736 Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov,
737 and Thomas Scialom. Llama 2: Open foundation and fine-
738 tuned chat models, 2023. 3
- 739 [45] Michael T Ullman. A neurocognitive perspective on lan-
740 guage: The declarative/procedural model. *Nature reviews*
741 *neuroscience*, 2(10):717–726, 2001. 4
- 742 [46] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai,
743 Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and
744 Hongxia Yang. Ofa: Unifying architectures, tasks, and
745 modalities through a simple sequence-to-sequence learning
746 framework. In *International Conference on Machine Learn-*
747 *ing*, pages 23318–23340. PMLR, 2022. 3
- 748 [47] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhil-
749 iang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mo-
750 hammed, Saksham Singhal, Subhajit Som, et al. Image as a
751 foreign language: Beit pretraining for all vision and vision-
752 language tasks. *arXiv preprint arXiv:2208.10442*, 2022. 3
- 753 [48] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang,
754 Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie
755 Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One
756 single transformer to unify multimodal understanding and
757 generation. *arXiv preprint arXiv:2408.12528*, 2024. 6, 7,
758 8
- 759 [49] Xiangyuan Xue, Zeyu Lu, Di Huang, Wanli Ouyang, and
760 Lei Bai. Genagent: Build collaborative ai systems with auto-
761 mated workflow generation—case studies on comfyui. *arXiv*
762 *preprint arXiv:2409.01392*, 2024. 2, 3, 4, 6, 7, 8
- 763 [50] Rui Yan, Mike Zheng Shou, Yixiao Ge, Alex Jinpeng
764 Wang, Xudong Lin, Guanyu Cai, and Jinhui Tang. Video-
765 text pre-training with learned regions. *arXiv preprint*
766 *arXiv:2112.01194*, 2021. 3
- 767 [51] Jinyu Yang, Jiali Duan, Son Tran, Yi Xu, Sampath Chanda,
768 Liqun Chen, Belinda Zeng, Trishul Chilimbi, and Junzhou
769 Huang. Vision-language pre-training with triple contrastive
770 learning. In *Proceedings of the IEEE/CVF Conference on*
771 *Computer Vision and Pattern Recognition*, pages 15671–
772 15680, 2022.
- 773 [52] Rowan Zellers, Jiasen Lu, Ximing Lu, Youngjae Yu, Yan-
774 peng Zhao, Mohammadreza Salehi, Aditya Kusupati, Jack
775 Hessel, Ali Farhadi, and Yejin Choi. Merlot reserve: Neu-
776 ral script knowledge through vision and language and sound.
777 In *Proceedings of the IEEE/CVF Conference on Computer*
778 *Vision and Pattern Recognition*, pages 16375–16387, 2022.
- 779 [53] Yan Zeng, Xinsong Zhang, and Hang Li. Multi-Grained
780 Vision Language Pre-Training: Aligning Texts with Visual
781 Concepts. In *International Conference on Machine Learn-*
782 *ing*, pages 25994–26009. PMLR, 2022. 3
- 783 [54] Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong
784 Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang,
785 Linyang Li, Hang Yan, Jie Fu, Tao Gui, Tianxiang Sun, Yu-
786 gang Jiang, and Xipeng Qiu. Anygpt: Unified multimodal
787 llm with discrete sequence modeling, 2024. 2, 3
- 788 [55] Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu
789 Wang, Yaqian Zhou, and Xipeng Qiu. Speechgpt: Empow-
790 ering large language models with intrinsic cross-modal con-
791 versational abilities, 2023. 3

Supplementary Material: Symbolic Representation for Any-to-Any Generative Tasks

Anonymous CVPR submission

Paper ID 3011

In Appendix A, we provide comprehensive details regarding the generative tasks in Sec. A.1 and present a thorough elaboration of the user study methodology in Sec. A.2. In Appendix B, we conduct an extensive qualitative analysis comparing our approach with existing unified multimodal frameworks, focusing on the quality of generated results. In Appendix C, we present additional experimental investigations, including a detailed comparison of computational efficiency versus human expert evaluation in Sec. C.1, and an in-depth analysis of examples and specific effects across three distinct syntax options in Sec. C.2. Finally, we examine the broader societal implications in Appendix D and discuss current limitations along with future research directions in Appendix E. Code and dataset are available at Anonymous Repository.

A. Experimental setup

A.1. Details on generative tasks

Our framework was evaluated across 12 distinct generative tasks collected from ComfyUI Examples [1], as detailed in Table 1. For Image2Mesh task, input images for mesh generation tasks were obtained from ComfyUI-3D-Pack [9], whereas other images were sourced from public repositories such as Vecteezy, Pexels, and Freepik. The evaluated tasks encompass a range of transformations, including:

- **Inpainting:** The task of image inpainting involves filling in appropriate content in the erased regions of a given image to generate a complete and visually coherent output.
- **Outpainting:** The image outpainting task extends the given image by generating a larger scene that seamlessly extends beyond the original boundaries while maintaining visual consistency.
- **Novel View Synthesis:** Novel view synthesis task takes a single object image as input and generates images of the object from novel viewpoints by inferring 3D geometric relationships from the 2D input.
- **Image merge:** The image merge task combines two landscape images to generate a new image that inherits the visual characteristics and features from both input images

harmoniously.

- **Merge model:** Merge model task blends different checkpoints for text-to-image generation models, enabling the creation of images that exhibit a combination of diverse visual styles and features.
- **Image2Mesh:** Image2Mesh task involves creating a 3D mesh model that corresponds to the given input image, capturing its geometric structure.
- **Multi-view image:** Given a single image, the multi-view task produces images of the same object from multiple viewpoints, offering comprehensive visual perspectives.
- **Image2Video:** Image2Video task creates a video sequence that is semantically related to the input image, expanding the static visual content into a dynamic narrative.
- **Text2Audio:** Text2Audio task enables the creation of music with specific styles based on the atmosphere and emotions conveyed in the textual description, facilitating text-guided music composition.
- **Text2Image:** Text2Image task synthesizes high-resolution, photorealistic images with rich details and cinematic quality based on the provided textual descriptions.
- **Text2Mesh:** Text2Mesh task creates 3D model mesh files that correspond to the given textual descriptions, translating language into three-dimensional geometric representations.
- **Text2Video:** Text2Video task involves creating video clips that align with the content and narrative described in the given text, bringing the written concepts to life through moving visuals.

A.2. User study setup

We selected five evaluators from diverse academic and cultural backgrounds¹. To ensure objectivity, we employed a double-blind evaluation method, ensuring that the evaluators were unaware of the model source for each result and that the presentation order was randomized. For each generative task, we conducted a one-on-one user study comparing

¹All evaluators were compensated with a wage of at least 30 US dollars per hour, which is higher than the statutory rate.

Table 1. **Task description.** Each line includes representative task examples and input-output modality pair for the task.

ID	Category	Example of natural language instruction	Input type	Output type
1	Inpainting	You are given an image named ‘yosemite_inpaint_example.png’. This image has had part of it erased, please inpaint a woman at the erased part to output a complete image called ‘woman_inpainted’.	Image	Image
2	Outpainting	You are given a image named ‘yosemite_outpaint_example.png’. Please outpaint the scenery of the given image and output a image called ‘scene_outpainted’.	Image	Image
3	Image merge	Given two images, ‘mountains.png’ and ‘sunset.png’, extract their visual features. Then, combine the extracted visual features to generate an image that depicts a beautiful scene with features from both input images. Finally, save the generated image as a file named ‘BeautifulScene’.	Image	Image
4	Novel view synthesis	You are given an image named ‘marble_statue.jpg’, please generate an image of the same object but from a different point of view. Save the output image as ‘statue_different_view’	Image	Image
5	Merge model	Generate an image with bottles containing a galaxy-like visual effect. Please merge two different checkpoints. Save the generated image as ‘galaxy_bottles’.	Text	Image
6	Image2Mesh	You are given an image named ‘marble_statue.jpg’. Please generate its 3D mesh and save the mesh as ‘marble_statue_mesh.obj’	Image	Mesh
7	Multi-view image	You are given an image named ‘marble_statue_rgba.png’. Please generate its multi-view images. The generated images‘filename prefix should be ‘Comfyui’.	Image	Image
8	Image2Video	You are given an image named ‘mountains.png’. Please create a 14 frame video of beautiful scenery from it.	Image	Video
9	Text2Audio	Generate an electronic dance music audio file inspired by a theme of ‘heaven church.’ Use an empty latent audio sample as the base, apply conditioning from a text description, and finally save the generated audio file as ‘electronic_audio’.	Text	Audio
10	Text2Image	Generate a high-resolution, cinematic image of an anthropomorphic fox in a sci-fi spaceship, wearing a spacesuit, with dramatic lighting and detailed features. The style should be realistic, high quality, in 4k resolution.	Text	Image
11	Text2Mesh	Generate a 3D mesh of a anime girl with short skirt and daisy blue eyes and save the mesh as ‘cute_girl.obj’.	Text	Mesh
12	Text2Video	Create a video of a cup of coffee being poured, but instead of coffee, a miniature galaxy swirls out, with stars and planets floating in the liquid.	Text	Video

our framework with all state-of-the-art unified multimodal frameworks (Show-o [12], SEED-X [3], LWM [7], Unified-IO 2 [8]). Detailed evaluation guidelines were provided, incorporating two ranking criteria. Evaluators assigned ranks from 1 (best) to n under each criterion.

- ① For **text-result alignment** assessment, evaluators were asked to assess if generated results faithfully represented all key elements specified in the input instructions (including scene composition, objects, styles and effects).

The evaluators should consider whether any required elements were missing or if there were unintended additions.

- ② For **result quality** assessment, evaluators evaluated the overall visual quality independent of the instructions. They focused on technical aspects like image sharpness, consistency of style, composition balance, and level of detail. For video outputs, they also considered motion smoothness and temporal coherence.

To evaluate performance uniformly, we averaged the rank-

085

086

087

088

089

090

091

092

093

Table 2. Average time cost (in seconds) compared to human ComfyUI experts. “NVS” denotes “Novel View Synthesis”, “I-2-3D” denotes “Image to 3D Mesh”, “T2M” denotes “Text to Mesh”.

Method	Inpaint	Outpaint	Img merge	NVS	Merge model	I-2-3D
Human	497.25	466.50	773.00	646.00	737.67	-
Ours	62.00	29.60	79.70	37.50	40.80	117.30
Speed up	8.02×	15.76×	9.70×	17.23×	18.08×	-
Method	T2I	T2A	Multi-view img	I2V	T2M	T2V
Human	537.25	278.00	-	590.00	-	1065.00
Ours	36.10	41.80	43.40	116.50	155.60	72.00
Speed up	14.88×	6.65×	-	5.06×	-	14.79×

ings for result quality and text-result alignment, with lower scores indicating better performance, with 1 being the best.

B. Qualitative comparison

In Figure 2 to 13, we compare our method with mainstream unified models (Show-o [12], SEED-X [3], LWM [7], Unified-IO 2 [8], i-Code-V3 [11] and AnyGPT [14]) across different generation tasks. Due to our model’s compositional nature, the LMs can invoke the most specialized functions and configure optimal parameters for each task. This flexibility in \mathcal{A} -language’s functions and parameters enables superior task-specific adaptation compared to implicit neural representations, which use identical weights and frameworks across all tasks. Furthermore, our framework eliminates the need for multi-task trade-offs in design, resulting in performance levels comparable to single-task expert models - all achieved in a training-free setting.

C. Additional experiments

C.1. Time cost: Human experts vs. Ours

Setup To evaluate the efficiency of our proposed method, four human experts, proficient in ComfyUI workflow construction, were invited to build the 12 generative workflows from scratch. We conducted experiments comparing the average time cost (in seconds) required by these human experts and our method for the corresponding tasks.²

- ① The timing started from the moment the experts saw the task and ended when they produced a result that matched the instructions.
- ② We provided reference key functions and parameters, and allowed the experts to use any online tools. They were not allowed to directly copy existing workflows to the workspace, but had to construct their own workflows based on the reference information.
- ③ The human experts were not required to optimize the workflows or improve the quality of the outcomes, they

²Since both human experts and the proposed method are based on the ComfyUI platform, their achieved results are not significantly different in quality. Therefore, it is not necessary to compare the quality differences.

Instruction: Given two images, ‘beach.png’ and ‘space_nebula.png’, extract their visual features. Then, combine the extracted visual features to generate an image of a beach with a surreal nebula sky. Finally, save the generated image as a file named ‘SpaceBeach.png’.



Figure 1. Example of image merge task.

only needed to complete the tasks.

Results As shown in Table 2, compared to the human experts with over 1 years of experience, our method achieved an efficiency improvement of **5-18 times**.

- ① **Source of efficiency gains:** Human experts, regardless of their years of experience, inevitably require substantial time for task contemplation and workspace manipulation. This inherent time investment cannot be eliminated from their workflow. In contrast, our approach leverages pre-trained LMs to generate workflows instantaneously, completing the task within seconds and eliminating the cognitive overhead and manual ComfyUI interface operations. The primary time expenditure is workflow execution and subsequent refinement phases.
- ② **Fluctuations in time costs:** Human experts exhibit significant variations in task completion times, ranging from 278.00 seconds for Text2Audio to 1065.00 seconds for Text2Video tasks, primarily due to the varying complexity of problem-solving and debugging processes. In contrast, our LM-based inference approach maintains nearly constant cognitive processing time, with time variations primarily attributed to workflow execution and refinement phases. This results in substantially smaller fluctuations, spanning from 29.60 seconds for Outpainting to 155.00 seconds for Text2Mesh tasks.

C.2. Details comparison with different syntaxes

To illustrate the differences between the syntaxes of \mathcal{A} -Language, we use the Image merge task as an example. The input, output, and instruction for this task can be found in Figure 1, while examples of the three different syntaxes are presented in Sections C.2.1 to C.2.3.

In Table 3, we compare three syntax styles for \mathcal{A} -Language. The Declarative syntax demonstrates superior overall performance, achieving an average compile rate of 0.97 and an execute rate of 0.77.

163

C.2.1. Example of dataflow syntax

```

# create nodes by instantiation
clipvisionencode_13 = CLIPVisionEncode()
clipvisionencode_36 = CLIPVisionEncode()
emptylatentimage_5 = EmptyLatentImage(width=768, height=768, batch_size=1)
unclipcheckpointloader_32 = unCLIPCheckpointLoader(ckpt_name='sd21-unclip-h.ckpt')
ksampler_3 = KSampler(seed=947446491266673, control_after_generate='randomize', steps=26, cfg=8,
↪ sampler_name='uni_pc_bh2', scheduler='normal', denoise=1)
cliptextencode_6 = CLIPTextEncode(text='beach with a surreal nebula sky')
cliptextencode_7 = CLIPTextEncode(text='boring, drab')
unclipconditioning_19 = unCLIPConditioning(strength=0.5, noise_augmentation=0.4)
unclipconditioning_37 = unCLIPConditioning(strength=0.5, noise_augmentation=0.4)
loadimage_beach = LoadImage(image='beach.png')
loadimage_nebula = LoadImage(image='space_nebula.png')
vaedecode_8 = VAE Decode()
saveimage_result = SaveImage(filename_prefix='SpaceBeach')

# link nodes by invocation
model_32, clip_32, vae_32, clip_vision_32, name_string_32 = unclipcheckpointloader_32()
image_beach, mask_beach = loadimage_beach()
image_nebula, mask_nebula = loadimage_nebula()
clip_vision_output_13 = clipvisionencode_13(clip_vision=clip_vision_32, image=image_beach)
clip_vision_output_36 = clipvisionencode_36(clip_vision=clip_vision_32, image=image_nebula)
conditioning_6 = cliptextencode_6(clip=clip_32)
negative_conditioning_7 = cliptextencode_7(clip=clip_32)
conditioning_19 = unclipconditioning_19(conditioning=conditioning_6,
↪ clip_vision_output=clip_vision_output_13)
conditioning_37 = unclipconditioning_37(conditioning=conditioning_19,
↪ clip_vision_output=clip_vision_output_36)
latent_5 = emptylatentimage_5()
latent_3 = ksampler_3(model=model_32, positive=conditioning_37, negative=negative_conditioning_7,
↪ latent_image=latent_5)
image_8 = vae_decode_8(samples=latent_3, vae=vae_32)
result = saveimage_result(images=image_8)

```

164

C.2.2. Example of declarative syntax

```

# Add Node
workflow.add_node("clipvisionencode_13", "CLIPVisionEncode", {})
workflow.add_node("emptylatentimage_5", "EmptyLatentImage", {"width": 768, "height": 768,
↪ "batch_size": 1})
workflow.add_node("unclipcheckpointloader_32", "unCLIPCheckpointLoader", {"ckpt_name":
↪ 'sd21-unclip-h.ckpt'})
workflow.add_node("clipvisionencode_36", "CLIPVisionEncode", {})
workflow.add_node("ksampler_3", "KSampler", {"seed": 947446491266673, "control_after_generate":
↪ 'randomize', "steps": 26, "cfg": 8, "sampler_name": 'uni_pc_bh2', "scheduler": 'normal',
↪ "denoise": 1})
workflow.add_node("cliptextencode_6", "CLIPTextEncode", {"text": 'beautiful photograph'})
workflow.add_node("cliptextencode_7", "CLIPTextEncode", {"text": 'bad hands'})
workflow.add_node("unclipconditioning_19", "unCLIPConditioning", {"strength": 0.5,
↪ "noise_augmentation": 0.4000000000000002})
workflow.add_node("unclipconditioning_37", "unCLIPConditioning", {"strength": 0.5,
↪ "noise_augmentation": 0.4000000000000002})
workflow.add_node("loadimage_1", "LoadImage", {"image": 'beach.png'})
workflow.add_node("loadimage_2", "LoadImage", {"image": 'space_nebula.png'})
workflow.add_node("vae_decode_8", "VAE Decode", {})
workflow.add_node("saveimage_9", "SaveImage", {"filename_prefix": 'SpaceBeach'})

```

```
# Invoke Node
workflow.invoke_node(["image_1", "mask_1"], "loadimage_1")
workflow.invoke_node(["image_2", "mask_2"], "loadimage_2")
workflow.invoke_node(["model_32", "clip_32", "vae_32", "clip_vision_32", "name_string_32"],
    ↪ "unclipcheckpointloader_32")
workflow.invoke_node(["latent_5"], "emptylatentimage_5")
workflow.invoke_node(["clip_vision_output_13"], "clipvisionencode_13")
workflow.invoke_node(["clip_vision_output_36"], "clipvisionencode_36")
workflow.invoke_node(["conditioning_6"], "cliptextencode_6")
workflow.invoke_node(["conditioning_7"], "cliptextencode_7")
workflow.invoke_node(["conditioning_19"], "unclipconditioning_19")
workflow.invoke_node(["conditioning_37"], "unclipconditioning_37")
workflow.invoke_node(["latent_3"], "ksampler_3")
workflow.invoke_node(["image_8"], "vaedecode_8")

# Link Node
workflow.connect("clip_vision_32", "clipvisionencode_13", "clip_vision")
workflow.connect("image_1", "clipvisionencode_13", "image")
workflow.connect("clip_vision_32", "clipvisionencode_36", "clip_vision")
workflow.connect("image_2", "clipvisionencode_36", "image")
workflow.connect("clip_32", "cliptextencode_6", "clip")
workflow.connect("clip_32", "cliptextencode_7", "clip")
workflow.connect("conditioning_6", "unclipconditioning_19", "conditioning")
workflow.connect("clip_vision_output_13", "unclipconditioning_19", "clip_vision_output")
workflow.connect("conditioning_7", "unclipconditioning_37", "conditioning")
workflow.connect("clip_vision_output_36", "unclipconditioning_37", "clip_vision_output")
workflow.connect("conditioning_19", "ksampler_3", "positive")
workflow.connect("conditioning_37", "ksampler_3", "negative")
workflow.connect("model_32", "ksampler_3", "model")
workflow.connect("latent_5", "ksampler_3", "latent_image")
workflow.connect("latent_3", "vaedecode_8", "samples")
workflow.connect("vae_32", "vaedecode_8", "vae")
workflow.connect("image_8", "saveimage_9", "images")
```

C.2.3. Example of pseudo-natural syntax

165

```
# create nodes by instantiation
clipvisionencode_13 is CLIPVisionEncode()
clipvisionencode_36 is CLIPVisionEncode()
emptylatentimage_5 is EmptyLatentImage with the parameters of (width is 768, height is 768, batch_size
    ↪ is 1)
unclipcheckpointloader_32 is unCLIPCheckpointLoader with the parameters of (ckpt_name is
    ↪ 'sd21-unclip-h.ckpt')
ksampler_3 is KSampler with the parameters of (seed is 947446491266673, control_after_generate is
    ↪ 'randomize', steps is 26, cfg is 8, sampler_name is 'uni_pc_bh2', scheduler is 'normal', denoise
    ↪ is 1)
cliptextencode_6 is CLIPTextEncode with the parameters of (text is 'beach with a surreal nebula sky')
cliptextencode_7 is CLIPTextEncode with the parameters of (text is 'boring, drab')
unclipconditioning_19 is unCLIPConditioning with the parameters of (strength is 0.5,
    ↪ noise_augmentation is 0.4)
unclipconditioning_37 is unCLIPConditioning with the parameters of (strength is 0.5,
    ↪ noise_augmentation is 0.4)
loadimage_beach is LoadImage with the parameters of (image is 'beach.png')
loadimage_nebula is LoadImage with the parameters of (image is 'space_nebula.png')
vaedecode_8 is VAE Decode()
saveimage_result is SaveImage with the parameters of (filename_prefix is 'SpaceBeach')

# link nodes by invocation
model_32, clip_32, vae_32, clip_vision_32, name_string_32 is unclipcheckpointloader_32()
```

```
image_beach, mask_beach is loadimage_beach()
image_nebula, mask_nebula is loadimage_nebula()
clip_vision_output_13 is clipvisionencode_13 with the parameters of (clip_vision is clip_vision_32,
↪ image is image_beach)
clip_vision_output_36 is clipvisionencode_36 with the parameters of (clip_vision is clip_vision_32,
↪ image is image_nebula)
conditioning_6 is cliptextencode_6 with the parameters of (clip is clip_32)
negative_conditioning_7 is cliptextencode_7 with the parameters of (clip is clip_32)
conditioning_19 is unclipconditioning_19 with the parameters of (conditioning is conditioning_6,
↪ clip_vision_output is clip_vision_output_13)
conditioning_37 is unclipconditioning_37 with the parameters of (conditioning is conditioning_19,
↪ clip_vision_output is clip_vision_output_36)
latent_5 is emptylatentimage_5()
latent_3 is ksampler_3 with the parameters of (model is model_32, positive is conditioning_37,
↪ negative is negative_conditioning_7, latent_image is latent_5)
image_8 is vaedecode_8 with the parameters of (samples is latent_3, vae is vae_32)
result is saveimage_result with the parameters of (images is image_8)
```

166

D. Social impacts

167

While mainstream research focuses on larger single-weight models on leaderboards, real-world AI application practices [2, 10], particularly in the multimodal content generation domain, increasingly employ composite workflows with multiple components, especially those based on platforms like ComfyUI and Blender. This paper distills the essential elements of these composite workflows into a simple symbolic language, allowing pre-trained language models to directly infer the workflows. This modular approach enables developers to create AIGC workflows with minimal coding effort.

178

However, this approach has potential negative impacts on the AIGC field. Firstly, it could change the production methods of AIGC practitioners, leading to a shift in the modes of labor within the domain. Traditional AIGC workflows often require practitioners to manually combine and adjust individual components, whereas using symbolic languages and pre-trained models to infer workflows can automate this process, reducing reliance on manual operations. Secondly, as the degree of automation increases, the amount of labor required in the field may decrease, impacting employment to a certain extent.

189

Despite these concerns, adopting symbolic languages and pre-trained models to infer workflows still has significant advantages. It can lower the entry barrier for AIGC workflow development, allowing more people to participate in the field. Moreover, by automating some repetitive and time-consuming work, practitioners can focus more on creativity and optimization, improving production efficiency and content quality.

197

E. Limitation and future work

198

While this paper establishes the foundational concepts and definitions for the formal language and inference method, further research and development are needed to bridge the gap between the proposed approach and real-world applications. This may involve incorporating domain-specific knowledge, integrating with hierarchical designs [5], tree search techniques [6], and advanced agentic learning-based strategies [15], and addressing practical concerns such as computational efficiency, user experience, and system robustness.

208

References

209

- [1] ComfyAnonymous. Comfyui examples. https://comfyanonymous.github.io/ComfyUI_examples/, 2023.
- [2] Yilun Du and Leslie Kaelbling. Compositional generative modeling: A single model is not all you need. *arXiv preprint arXiv:2402.01103*, 2024.
- [3] Yuying Ge, Sijie Zhao, Jinguo Zhu, Yixiao Ge, Kun Yi, Lin Song, Chen Li, Xiaohan Ding, and Ying Shan. Seed-x: Mul-

timodal models with unified multi-granularity comprehension and generation. *arXiv preprint arXiv:2404.14396*, 2024.

[4] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *CVPR*, pages 14953–14962, 2023.

[5] Sirui Hong, Yizhang Lin, Bangbang Liu, Biniao Wu, Danyang Li, Jiaqi Chen, Jiayi Zhang, Jinlin Wang, Lingyao Zhang, Mingchen Zhuge, et al. Data interpreter: An llm agent for data science. *arXiv preprint arXiv:2402.18679*, 2024.

[6] Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. Tree search for language model agents. *arXiv preprint arXiv:2407.01476*, 2024.

[7] Hao Liu, Wilson Yan, Matei Zaharia, and Pieter Abbeel. World model on million-length video and language with ringattention. *arXiv preprint*, 2024.

[8] Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision, language, audio, and action, 2023.

[9] MrForExample. Comfyui-3d-pack. <https://github.com/MrForExample/ComfyUI-3D-Pack>, 2024.

[10] Lin Qiao. Fireworks ai raises 52m series b to lead industry shift to compound ai systems. <https://fireworks.ai/blog/fireworks-ai-series-b-compound-ai>, 2024.

[11] Zineng Tang, Ziyi Yang, Chenguang Zhu, Michael Zeng, and Mohit Bansal. Any-to-any generation via composable diffusion. *arXiv preprint arXiv:2305.11846*, 2023.

[12] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024.

[13] Xiangyuan Xue, Zeyu Lu, Di Huang, Wanli Ouyang, and Lei Bai. Genagent: Build collaborative ai systems with automated workflow generation—case studies on comfyui. *arXiv preprint arXiv:2409.01392*, 2024.

[14] Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang, Linyang Li, Hang Yan, Jie Fu, Tao Gui, Tianxiang Sun, Yungang Jiang, and Xipeng Qiu. Anygpt: Unified multimodal llm with discrete sequence modeling, 2024.

[15] Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: Llm agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 19632–19642, 2024.

Table 3. **Performance comparison on syntax style.** We report the pass rate for a single run (Pass@1%). “Comp” denotes “Compile”, “Exec” denotes “Execute”.

Syntax	Inpainting		Img merge		Outpainting		Novel view syn.		Merge model		Img2Mesh		
	Comp	Exec	Comp	Exec	Comp	Exec	Comp	Exec	Comp	Exec	Comp	Exec	
Dataflow [4, 13]	1.0	0.9	1.0	0.7	0.9	0.9	0.7	0.5	1.0	0.9	1.0	1.0	
Pseudo-natural	0.9	0.9	0.9	0.8	0.9	0.9	0.8	0.7	0.9	0.9	0.8	0.7	
Declarative	1.0	0.7	0.9	0.7	1.0	1.0	0.9	0.8	0.9	0.9	1.0	0.8	
Multi-view img		Img2Video		Text2Audio		Text2Img		Text2Mesh		Text2Video		Average	
Comp	Exec	Comp	Exec	Comp	Exec	Comp	Exec	Comp	Exec	Comp	Exec	Comp	Exec
0.8	0.1	1.0	0.9	0.6	0.4	1.0	0.9	0.3	0.3	0.8	0.1	0.84	0.63
0.9	0.4	1.0	1.0	0.6	0.4	0.8	0.8	0.3	0.0	0.9	0.1	0.81	0.63
1.0	1.0	0.9	0.6	1.0	0.8	1.0	1.0	1.0	0.8	1.0	0.1	0.97	0.77

Table 4. **Natural language instructions used in merge model for image generation.** The “position” column indicates the image’s location in figure 9 grid using (row, column) coordinates, counting from top-left to bottom-right.

ID	Position	Natural Language Instruction
1	(1, 1)	Create a high-definition, futuristic image of a bustling neon-lit city at night, with towering skyscrapers, rain-soaked streets, and holographic billboards. The style should be cyberpunk, rich in vibrant colors and contrast. Please merge two different checkpoints.
2	(2, 1)	Produce an image of a Norse god standing atop a cliff with thunderous clouds and a glowing hammer. The scene should have dramatic, epic lighting in the style of classical oil paintings. Please merge two different checkpoints.
3	(3, 1)	Generate a highly detailed, 4K image of a steampunk inventor’s workshop, filled with intricate gears, brass machines, and soft, warm lighting. The style should be vintage and richly textured. Please merge two different checkpoints.
4	(4, 1)	Create a beautiful underwater scene featuring a bioluminescent jellyfish forest with mythical creatures swimming around. The image should have a mystical, tranquil feel with soft blue-green hues and glowing details. Please merge two different checkpoints.
5	(5, 1)	Design a minimalist, high-contrast image of a lone cactus in a vast desert under a giant, crimson sun. The colors should be bold, with a surreal, almost abstract aesthetic. Please merge two different checkpoints.
6	(1, 2)	Produce a hauntingly beautiful portrait of a Victorian woman in dark attire, surrounded by a foggy, candlelit room with antique furniture. The style should be Gothic, with a moody, mysterious vibe. Please merge two different checkpoints.
7	(2, 2)	Generate a detailed illustration of an animal tea party in a forest clearing, featuring animals like rabbits and foxes dressed in Victorian attire. The style should be whimsical, with soft, pastel colors and charming details. Please merge two different checkpoints.
8	(3, 2)	Create a high-resolution image of an alien planet landscape with two suns and strange rock formations, set against a sky filled with vibrant galaxies. The style should be sci-fi with vibrant colors and atmospheric lighting. Please merge two different checkpoints.
9	(4, 2)	Produce an image of a retro-futuristic city with flying cars and curved glass buildings, all in a 1980s-inspired color palette. The style should be bold, with neon hues and a sense of nostalgic futurism. Please merge two different checkpoints.

Generate a high-resolution, cinematic image of an anthropomorphic fox in a sci-fi spaceship, wearing a spacesuit, with dramatic lighting and detailed features. The style should be realistic, high quality, in 4k resolution.



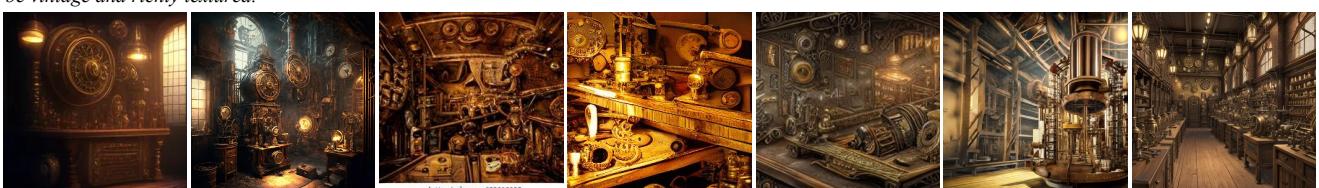
Create a high-definition, futuristic image of a bustling neon-lit city at night, with towering skyscrapers, rain-soaked streets, and holographic billboards. The style should be cyberpunk, rich in vibrant colors and contrast.



Produce an image of a Norse god standing atop a cliff with thunderous clouds and a glowing hammer. The scene should have dramatic, epic lighting in the style of classical oil paintings.



Generate a highly detailed, 4K image of a steampunk inventor's workshop, filled with intricate gears, brass machines, and soft, warm lighting. The style should be vintage and richly textured.



Create a beautiful underwater scene featuring a bioluminescent jellyfish forest with mythical creatures swimming around. The image should have a mystical, tranquil feel with soft blue-green hues and glowing details.



Show-o

SEED-X

LWM

Unified-IO 2

i-Code-V3

AnyGPT

Ours

Figure 2. Qualitative results of Text2Image task (Part 1).

Design a minimalist, high-contrast image of a lone cactus in a vast desert under a giant, crimson sun. The colors should be bold, with a surreal, almost abstract aesthetic.



Produce a hauntingly beautiful portrait of a Victorian woman in dark attire, surrounded by a foggy, candlelit room with antique furniture. The style should be Gothic, with a moody, mysterious vibe.



Generate a detailed illustration of an animal tea party in a forest clearing, featuring animals like rabbits and foxes dressed in Victorian attire. The style should be whimsical, with soft, pastel colors and charming details.



Create a high-resolution image of an alien planet landscape with two suns and strange rock formations, set against a sky filled with vibrant galaxies. The style should be sci-fi with vibrant colors and atmospheric lighting.



Produce an image of a retro-futuristic city with flying cars and curved glass buildings, all in a 1980s-inspired color palette. The style should be bold, with neon hues and a sense of nostalgic futurism.



Show-o

SEED-X

LWM

Unified-IO 2

i-Code-V3

AnyGPT

Ours

Figure 3. Qualitative results of Text2Image task (Part 2).

This image has had part of it erased, please inpainting a woman at the erased part to output a complete image.



Please inpaint a friendly alien standing beside the bench, and output a complete image.



Please inpaint a lush on the desk, and output a complete image.



Input

Show-o

Ours

Please inpaint a hat over on the main ancient statue, and output a complete image.



Please inpaint flowers with pink petals floating on the lake, and output a complete image.



Please inpaint a modern time traveler with a smartphone exploring the ruins, and output a complete image.



Please inpaint a rainbow-colored unicorn grazing in the meadow, and output a complete image.

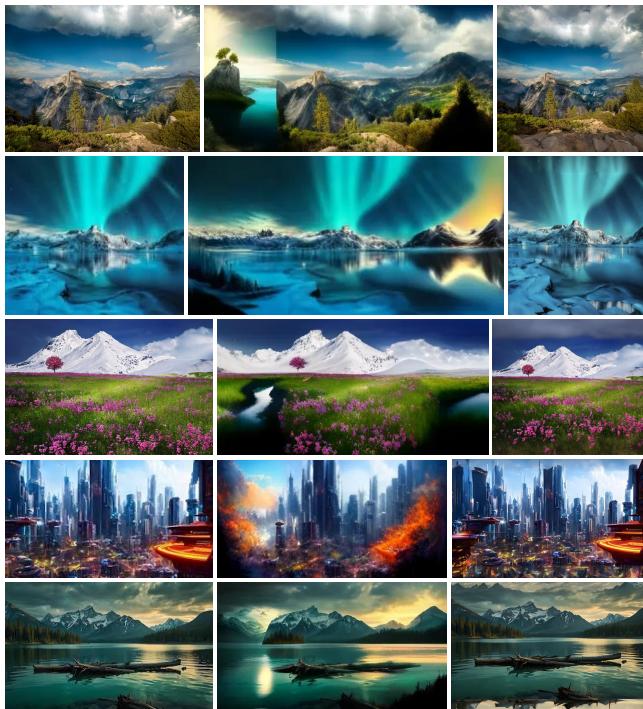


Input

Show-o

Ours

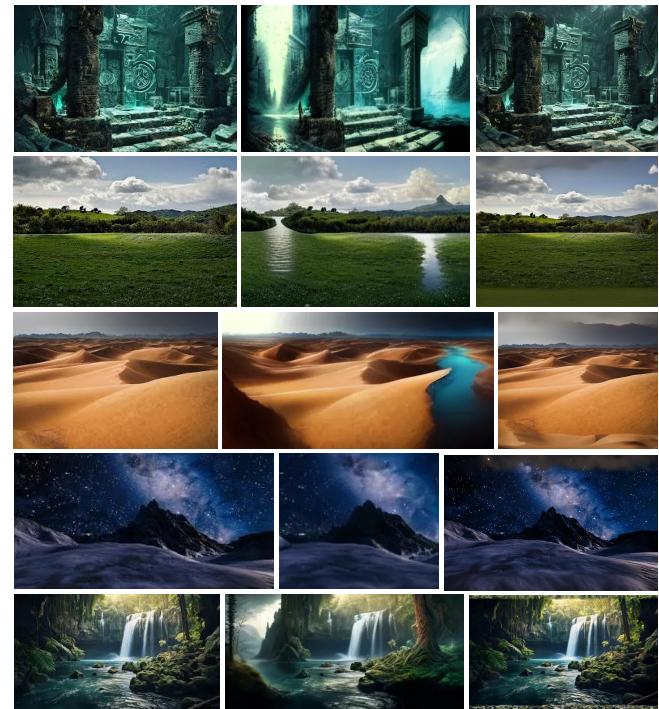
Figure 4. Qualitative results of image inpainting.



Input

Show-o

Ours



Input

Show-o

Ours

Figure 5. Qualitative results of image outpainting.



Figure 6. Qualitative results of image merge.



Figure 7. Qualitative results of Text2Video.

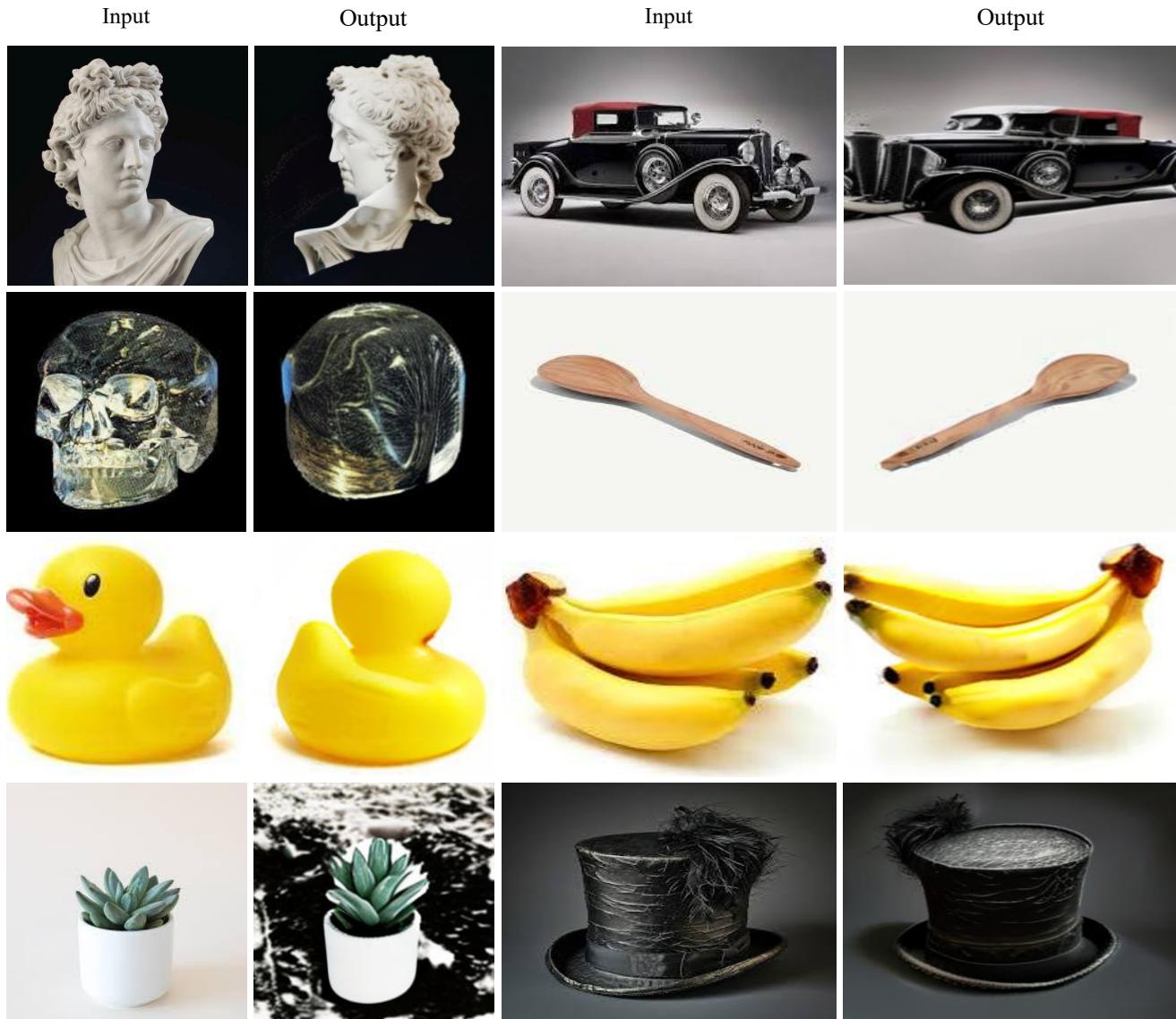


Figure 8. Qualitative results of novel view synthesis (NVS).



Figure 9. **Qualitative results of merge model.** This task allows visual style combinations through checkpoint blending. The visualization demonstrates generated outputs, where each image corresponds to its respective natural language instruction as detailed in Table 4.

Table 5. **Natural language instructions used in Text2Mesh generation.** The “position” column indicates the 3D mesh’s location in figure 13 grid using (row, column) coordinates, counting from top-left to bottom-right.

ID	Position	Natural Language Instruction
1	(1, 1)	Generate a 3D mesh of an anime girl with short skirt and daisy blue eyes and save the mesh as ‘cute_girl.obj’.
2	(2, 1)	Generate a 3D mesh of a plain ceramic coffee mug with a matte white finish. It features a gently curved, sturdy handle for gripping, a slightly rounded base, and a smooth, untextured surface that reflects faint ambient light. Save the mesh as ‘coffee_mug.obj’.
3	(3, 1)	Generate a 3D mesh of a classic hardcover book with a solid blue cover. The book has a subtle fabric texture and rounded corners. The pages are aligned neatly with a slight golden tint at the edges, giving a vintage look. Save the mesh as ‘blue_book.obj’.
4	(4, 1)	Generate a 3D mesh of a round, bright orange with a textured peel covered in small dimples. It has a tiny, dried green stem on top, and the surface shows a faint, shiny sheen, indicating juiciness. Save the mesh as ‘orange_fruit.obj’.
5	(5, 1)	Generate a 3D mesh of a simple black office chair with a flat seat and a low backrest. It has a minimalistic design, thin matte finish, and stands on a five-wheel base, with each wheel small and unobtrusive. Save the mesh as ‘office_chair.obj’.
6	(6, 1)	Generate a 3D mesh of a standard incandescent light bulb with a clear glass surface. The metallic base has grooved ridges for screwing in, and inside, a thin tungsten filament is suspended by two small metal wires. Save the mesh as ‘light_bulb.obj’.
7	(1, 2)	Generate a 3D mesh of a simple wooden spoon with a smooth, polished light brown surface. The handle is straight with a slight curve at the end, and the bowl of the spoon is shallow with a rounded edge. Save the mesh as ‘wooden_spoon.obj’.
8	(2, 2)	Generate a 3D mesh of a cylindrical transparent water bottle with a faint blue tint. It features a screw-on cap with ridges for grip, smooth body with slight indentations for holding, and tiny air bubbles trapped inside the water. Save the mesh as ‘water_bottle.obj’.
9	(3, 2)	Generate a 3D mesh of a small green apple with a shiny, waxy surface. It has a slightly irregular shape, a tiny brown stem, and a smooth skin with light speckles and green highlights. Save the mesh as ‘green_apple.obj’.
10	(4, 2)	Generate a 3D mesh of a traditional wooden pencil with a yellow hexagonal body. The pencil has a sharp graphite tip and a pink eraser on the other end, held by a shiny metal band. There are faint lines showing the wood grain pattern. Save the mesh as ‘yellow_pencil.obj’.

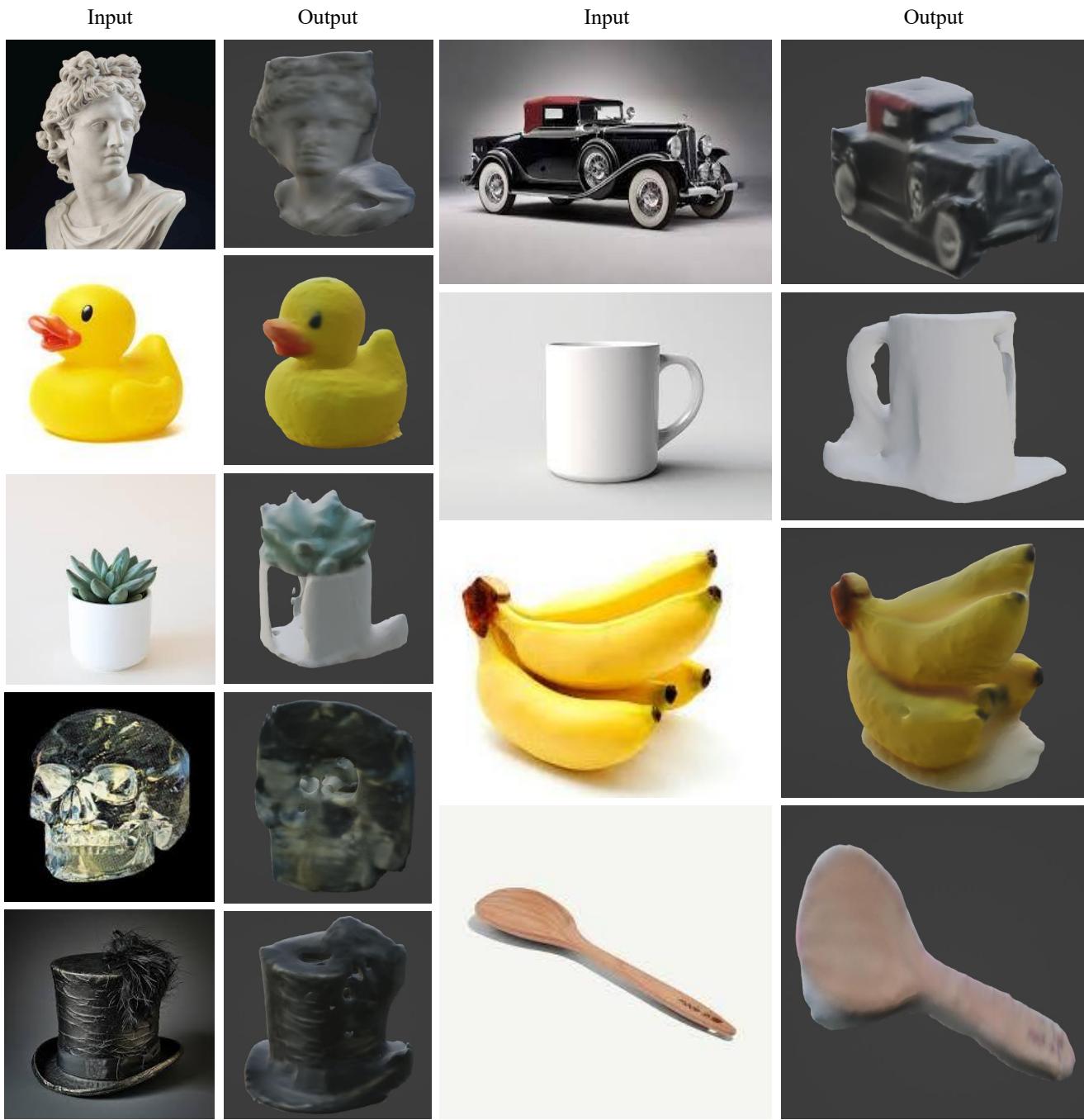


Figure 10. Qualitative results of Image2Mesh.



Figure 11. Qualitative results of multi-view image generation.

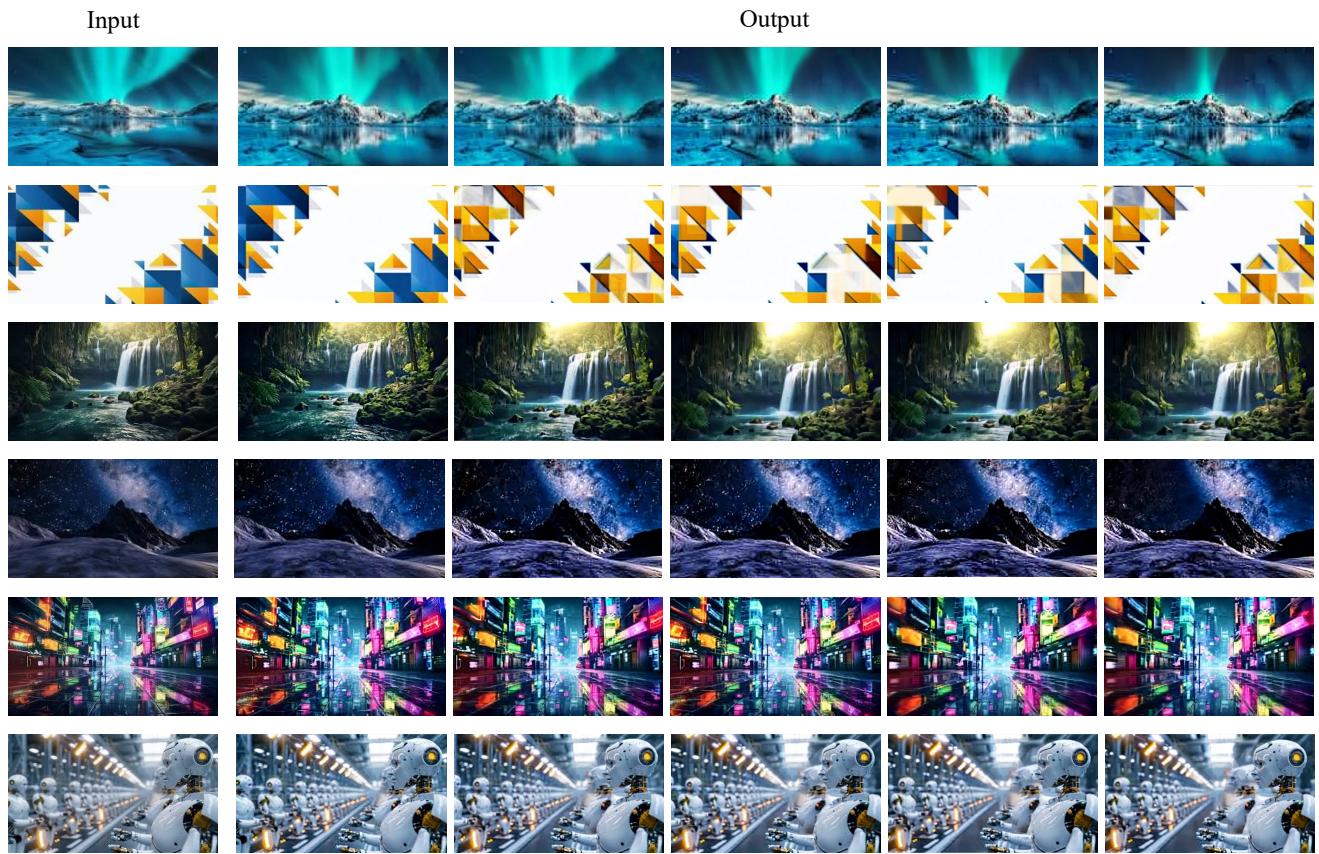


Figure 12. Qualitative results of Image2Video.

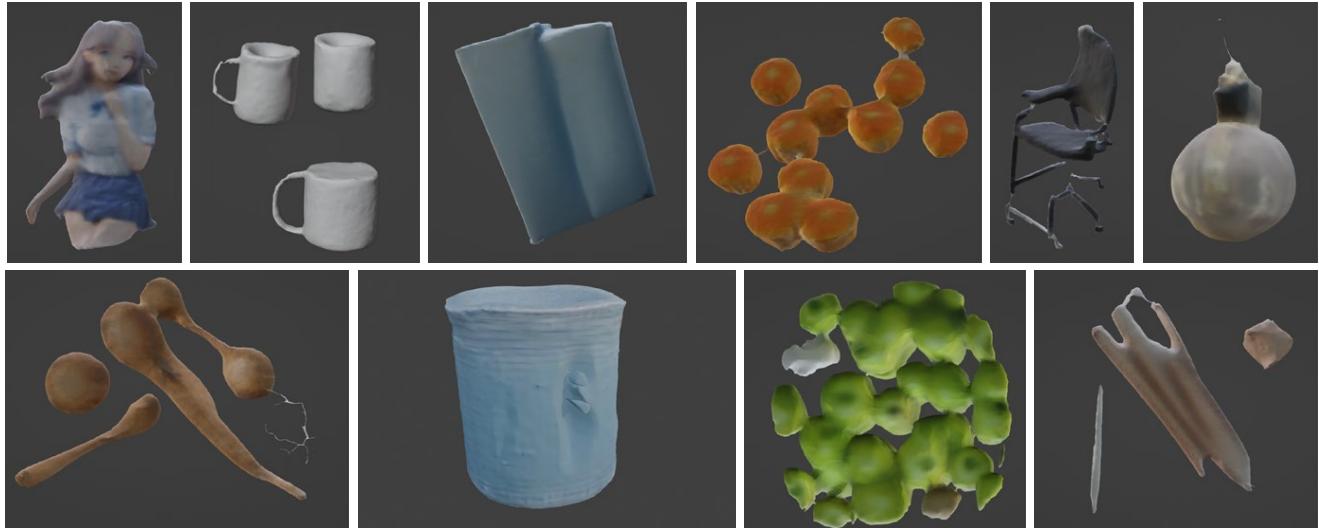


Figure 13. Qualitative results of Text2Mesh. The generated 3D meshes are synthesized based on their corresponding natural language instructions as specified in Table 5.