# Prompt-to-Product: Generative Assembly via Bimanual Manipulation

Ruixuan Liu*, Philip Huang*, Ava Pun, Kangle Deng, Shobhit Aggarwal, Kevin Tang, Michelle Liu,
Deva Ramanan, Jun-Yan Zhu, Jiaoyang Li and Changliu Liu

Carnegie Mellon University

*Abstract*—Creating assembly products demands significant manual effort and expert knowledge in 1) designing the assembly and 2) constructing the product. This paper introduces Prompt-to-Product, an automated pipeline that generates real-world assembly products from natural language prompts. Specifically, we leverage LEGO bricks as the assembly platform and automate the process of creating brick assembly structures. Given the user design requirements, Prompt-to-Product generates physically buildable brick designs, and then leverages a bimanual robotic system to construct the real assembly products, bringing user imaginations into the real world. We conduct a comprehensive user study, and the results demonstrate that Prompt-to-Product significantly lowers the barrier and reduces manual effort in creating assembly products from imaginative ideas.

## I. INTRODUCTION

Transforming design concepts into real products is a complex, time-consuming process that requires both creativity and technical expertise. Recent advances in generative artificial intelligence [26, 30] and additive manufacturing technologies [27, 1] have lowered the barrier to physical prototyping. While existing techniques predominantly focus on rigid objects with fully connected, monolithic structures, they fall short when applied to assembly objects, *i.e.*, products composed of multiple interlocking components. Unlike monolithic bodies, assembly objects play a crucial role in the real world, as most engineered products, from toys and furniture to machines and electronics, are inherently modular and require assembly. These objects cannot be replaced by single rigid bodies without sacrificing practicality or functionality. Bringing 3D assembly designs into physical reality is particularly challenging, as assembly products must not only meet visual or aesthetic expectations, but also satisfy strict *physical constraints*.

**Challenges.** As illustrated in Fig. 1, three important *physical constraints* need to be satisfied. First, the product needs to satisfy the *environmental resource constraint*. We must reason about the material available and build the assembly product only using the available inventory. Second, the product should satisfy the *embodiment dexterity constraint*. It is important to understand the capability of the system and know what can be physically built within the skill set. And most importantly, the assembly product should satisfy the *physical feasibility constraint*. It is critical to ensure that the assembly product is feasible as expected in the real world. Due to the physical constraints mentioned above, it is challenging to bring assembly design ideas to life as real products.

**Related Works** Recent work [17] generates 2D puzzles using available building blocks but requires human effort
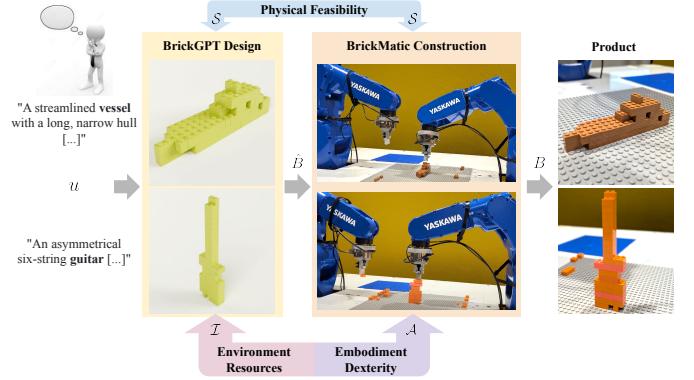
* Indicates equal contribution.



Fig. 1: **An Overview of Prompt-to-Product.** Given the user design ideas, input as text prompts, Prompt-to-Product 1) generates the assembly designs using BRICKGPT, and 2) physically constructs the assembly products with BRICKMATIC, bringing user imaginations into real products.

to build them. Zhang et al. [28] transform a monolithic 3D object into an assembly by decomposing it into sub-components. However, they assume an unlimited inventory and print [1] the components, which later require manual assembly. Our goal aligns closely with that of Kyaw et al. [11] and Goldberg et al. [3]. Kyaw et al. [11] generate an assembly design using prefabricated cuboctahedron blocks with magnet connections from user audio input, and build the assembly using a robot. However, using customized cuboctahedron blocks as the assembly platform limits the replicability. In addition, their system uses a simple rule-based check to assess physical stability and does not build collaboratively (*e.g.*, support), which prevents it from scaling to complex structures. On the other hand, Goldberg et al. [3] generate an assembly design from user text input and enable a robot to construct the structure physically using 3D printed blocks. However, these blocks have smooth surfaces without interconnections, and thus have limited expressiveness. Similarly, they use a single robot to construct the assembly object, which has very limited manipulation dexterity to build complex structures (*i.e.*, > 10 blocks). In this work, our goal is to develop an integrated system that creates real-world 3D assemblies from user text input, allowing easy reproducibility, high customizability, and scalability up to complex structures with up to hundreds of components.

**Problem Statement.** This paper studies the problem of translating natural language prompts into real-world 3D assembly products. In particular, we leverage the LEGO® brick system as the assembly platform. LEGO bricks are low-cost and readily available, allowing easy replication for benchmarking. In addition, they are inherently modular and highly reconfigurable, offering a vast design space through simple yet expressive

interconnections. This makes them uniquely well-suited for this study. By using LEGO bricks, this work aims to significantly reduce the manual effort required to turn user-generated text prompts into physical brick assembly structures—bridging the gap between imagination and realization. In the context of brick assembly, the physical constraints shown in Fig. 1 are defined as follows:

1) Environment resources: we constrain the available inventory to contain 8 types of widely used bricks, including $1 \times 1, 1 \times 2, 1 \times 4, 1 \times 6, 1 \times 8, 2 \times 2, 2 \times 4$, and $2 \times 6$.
2) Embodiment dexterity: we leverage general-purpose robotic arms with customized end-of-arm tools (EOAT) to build brick assemblies.
3) Physical feasibility: the brick assembly product should be physically stable and not collapse.

**Prompt-to-Product.** To this end, this paper introduces Prompt-to-Product, an automated pipeline that creates real-world brick assembly products from natural language prompts as illustrated in Fig. 1. Due to the inherent complexity of assembly tasks, end-to-end methods [8, 2, 6] often struggle with the long-horizon reasoning and fine-grained dexterity these tasks demand. Prompt-to-Product adopts a staged architecture to address the distinct yet interdependent challenges of translating natural language prompts into physical brick assemblies. This modular design separates the design generation and physical construction phases, allowing each to apply creative inference in the former and dexterous manipulation in the latter. Crucially, these modules are not isolated; they are tightly coupled through a shared physics reasoning module, which enforces global feasibility and maintains consistency across the pipeline.

In the first stage, BRICKGPT [22] generates brick assembly designs from user prompts, exploring the vast combinatorial space of brick structures while satisfying semantic and aesthetic user intents. In the second stage, BRICKMATIC physically realizes the designs using a bimanual robotic system capable of dexterously executing long-horizon manipulation tasks. At the core of Prompt-to-Product is a physics reasoning module, based on the structural stability analysis in [13], which plays a dual role. First, it constrains BRICKGPT to produce only physically buildable and stable designs given the available inventory. Second, it guides BRICKMATIC through the physical assembly process by reasoning over intermediate stability and ensuring that the product remains constructible at every step.

As illustrated in Fig. 1, this staged yet interconnected architecture enables Prompt-to-Product to effectively bridge user intent, structural design, and robotic execution—offering a coherent pipeline from prompt to product. We conduct a comprehensive user study, including more than 20 participants. The user study demonstrates that Prompt-to-Product significantly reduces the required manual effort and expert knowledge in creating brick assembly products from imaginative ideas.

**Contributions.** Our contributions are listed as follows:

1) We present Prompt-to-Product to create assembly products from users' design ideas in the form of text prompts.
2) We introduce BRICKMATIC, an integrated bimanual robotic system with enhanced dexterity capable of constructing customized brick assembly products.

3) We conduct a comprehensive user study with participants from different backgrounds. The results demonstrate that Prompt-to-Product effectively reduces manual effort in creating assembly products from abstract ideas.

## II. OVERVIEW OF PROMPT-TO-PRODUCT

We formulate the Prompt-to-Product system as a staged process that maps a high-level user text prompt into a physically realized brick assembly through a sequence of modules. The staged design reflects the fundamentally different reasoning skills required at each step: creative semantic generation, physical feasibility evaluation, sequential assembly planning, and parallel robotic collaboration. Importantly, the stages are tightly coupled to ensure global consistency and feasibility.

### A. Problem Formulation

We formulate the problem using the following notations:

- $\{\cdot\}$: a set of elements where the order does not matter.
- $[\cdot]$: a sequence of elements where the order matters.
- $u \in \mathcal{U}$: a user prompt in natural language describing the design requirements of the assembly product.
- $B$: a brick structure layout represented as a set of bricks:

$$B = \{b_1, b_2, \ldots, b_N\}$$

where $N$ is the number of bricks in the structure. Each brick $b_i = \{c_i, p_i, \omega_i\}$ where $c_i$ denotes the brick type, $p_i \in \mathbb{R}^3$ is the brick pose in space, and $\omega_i$ is the planar orientation of the brick.

- $\mathcal{I}$: the set of available inventory. The environmental resource constraint enforces $\forall i \in \{1, 2, \ldots, N\}, c_i \in \mathcal{I}$.
- $\mathcal{A}$: the set of skills that the system is capable of. The dexterity constraint requires the system to only use available skills $a_i \in \mathcal{A}$ to assemble $b_i$.
- $Q = [\{b_1^a, a_1\}, \{b_2^a, a_2\}, \ldots, \{b_N^a, a_N\}]$: the sequence to build a brick structure $B$. Each brick $b_i^a \in B$ and $a_i \in \mathcal{A}$ is the skill for assembling $b_i^a$.
- $G = \{V, E\}$: the bimanual system execution plan, represented as a temporal plan graph (TPG) [4]. Each node in $V$ denotes a robot action to execute the skill sequence $[a_1, a_2, \ldots, a_N]$, and each edge in $E$ indicates a precedence constraint between node executions.
- $S = [s_1; s_2; \ldots; s_N] \in \mathbb{R}^N$: the stability evaluation of the entire structure $B$. Each $s_i \in \mathbb{R}$ is the stability score of the brick $b_i$. A brick is non-collapsing if $s_i > 0$.
- $\mathcal{S} = \{B \mid \forall i \in \{1, 2, \ldots, N\}, s_i > 0\}$: the set of all physically stable brick designs. The physical feasibility constraint requires $B \in \mathcal{S}$.
- $\mathcal{B}_{\mathcal{I}, \mathcal{A}, \mathcal{S}}$: the space of all physically buildable designs with the given inventory and system (*i.e.*, structures that satisfy all physical constraints: $\mathcal{I}, \mathcal{A}, \mathcal{S}$).
- $\hat{\mathcal{B}}$: the set of all possible designs, either buildable or not.
- $\hat{\mathcal{B}}(u) \subseteq \hat{\mathcal{B}}$: the set of all brick assembly designs that are semantically valid given the prompt $u$.

Due to physical constraints, not all elements in $\hat{\mathcal{B}}(u)$ are constructible. We define the feasible subset:

$$\mathcal{B}(u \mid \mathcal{I}, \mathcal{A}, \mathcal{S}) \equiv \hat{\mathcal{B}}(u) \cap \mathcal{B}_{\mathcal{I}, \mathcal{A}, \mathcal{S}}$$

which includes only those assemblies that can be constructed with the given inventory and system capabilities.

The proposed Prompt-to-Product is a mapping from the user prompt $u$ to a physical assembly $B \in \mathcal{B}(u \mid \mathcal{I}, \mathcal{A}, \mathcal{S})$ that is buildable by the system with the given inventory. This paper focuses on constructing such a mapping efficiently, ensuring that the entire process, from prompt interpretation to physical realization, can be completed to support interactive and practical use cases. We also acknowledge that this is not the only way to instantiate such a mapping, but rather a structured and tractable approach enabled by modular reasoning.

### B. Physical Feasibility

Reasoning physical feasibility is crucial in both designing and constructing the brick assembly. It is important to ensure that the generated brick design is physically feasible in the first place; otherwise, it is impossible to build it. Similarly, it is also critical to understand the physical feasibility during construction so that the system can reliably build the brick assembly rather than collapse it. Evaluating the physical feasibility (stability) is non-trivial:

$$S = f_{\text{stability}}(B), \tag{1}$$

where $B \in \hat{\mathcal{B}}$ can be any brick structure. Conventional approaches [25, 3, 29] rely on physics engines to simulate the structural stability. However, existing simulations fail to accurately model the deformations and interconnections between bricks, and thus, cannot correctly estimate the structural stability. Pletz and Drvoderic [20] leverage the finite-element method to simulate high-fidelity behavior of brick structures, which is difficult to scale. Luo et al. [16] estimate the stability by analyzing the internal force distribution. However, their formulation only applies to structures that are connected as one piece, which is a subset of $\hat{\mathcal{B}}$. We leverage the stability analysis in [13] to model $f_{\text{stability}}(\cdot)$ and evaluate $S$ for any brick structure $B \in \hat{\mathcal{B}}$. The details of $f_{\text{stability}}(\cdot)$ are discussed in Sec. III.

### C. Stage 1: Prompt-to-Design via BRICKGPT

Given the input user prompt $u$, we define the generative mapping to generate the candidate brick layout design as:

$$\hat{B} = f_{\text{design}}(u),$$

where $\hat{B} = \{\hat{b}_1, \hat{b}_2, \ldots, \hat{b}_N\}$. The key requirement is that $\hat{B}$ should lie in $\mathcal{B}(u \mid \mathcal{I}, \mathcal{A}, \mathcal{S})$ in order to be physically built. Since $\mathcal{A}$ is highly system dependent, we relax the constraint and require $\hat{B} \in \mathcal{B}(u \mid \mathcal{I}, \mathcal{S})$ in this stage, and consider $\mathcal{A}$ in the second stage. By relaxing $\mathcal{A}$, the generated $\hat{B}$ becomes a human-buildable brick assembly design. Conventional approaches decouple the constraints, *i.e.*, they leverage generative models [26, 30] to first generate the 3D shape that is visually and semantically appealing given $u$, and then use mesh-to-brick algorithms [16, 15] to transform the 3D shape to a brick assembly structure that satisfy the physical constraints, *i.e.*, $\mathcal{I}$ and $\mathcal{S}$. However, the generated 3D shape may have no solution for mesh-to-brick, meaning that it cannot be built using the brick inventory while being stable. Hence, we leverage
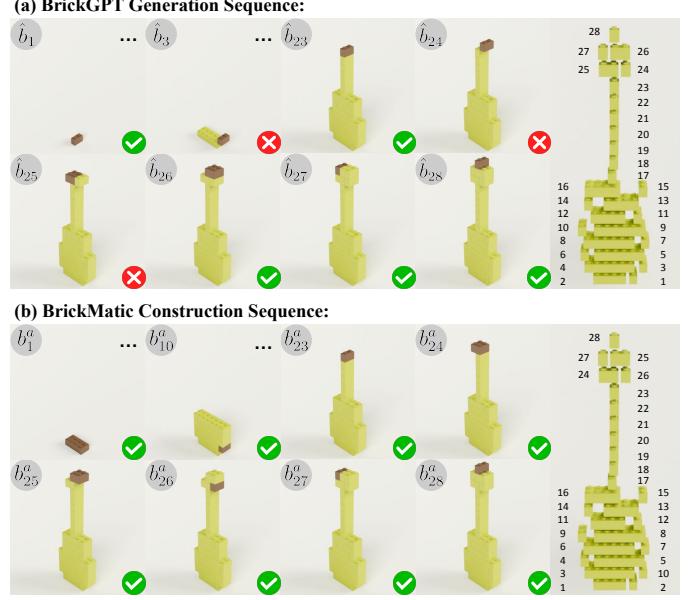


Fig. 2: **Brick Order.** (a) The sequence of BRICKGPT generating the design $\hat{B}$. (b) The assembly sequence $Q$ planned by BRICKMATIC to physically construct the structure. The check mark indicates that the partial structure satisfies the physical feasibility constraint.

BRICKGPT [22], an end-to-end approach, to generate the brick design $\hat{B} \in \mathcal{B}(u \mid \mathcal{I}, \mathcal{S})$ as

$$\begin{aligned} \hat{B} &= f_{\text{BRICKGPT}}(u), \\ \text{s.t.} \quad & c_i \in \mathcal{I}, \ \forall i \in \{1, 2, \ldots, N\}, \\ & \hat{B} \in \mathcal{S}, \end{aligned} \tag{2}$$

where $N$ is internally determined by $f_{\text{BRICKGPT}}(\cdot)$. The method for $f_{\text{BRICKGPT}}(\cdot)$ is discussed in Sec. IV.

### D. Stage 2: Design-to-Product via BRICKMATIC

Given the design $\hat{B}$, the builder system constructs the brick assembly product as

$$B = f_{\text{build}}(\hat{B})$$

At this step, the brick layout is not to be changed (essentially $B = \hat{B}$, though one is in the real physical space and the other is in the virtual space). But the building procedure needs to be carefully reasoned, taking the consideration of the physical constraints. Since the resource constraint is enforced in the first stage, we relax the constraint and only account for constraints $\mathcal{A}$ and $\mathcal{S}$. Existing robotic systems often fall short with long-horizon reasoning and fine-grained dexterity, and thus, they only perform in simulation [21, 18] or build simple structures [14, 3]. Hence, we introduce BRICKMATIC, a bimanual robotic system capable of building customized brick assemblies. In particular, we expand $\mathcal{A}$ through embodiment design and skill learning discussed in Sec. V-A and V-B. The physical construction is accomplished by a multi-level reasoning framework. The first step is to reason over the assembly sequence $Q$, which enforces that each partial assembly $B_i^a = \{b_1^a, b_2^a, \ldots, b_i^a\} \in \mathcal{B}_{\mathcal{I}, \mathcal{A}, \mathcal{S}}$ is physically buildable given the system's capabilities. As shown in Fig. 2, the planned $b_i^a \in Q$ can be different from $\hat{b}_i \in \hat{B}$ resulted from BRICKGPT in Eq. (2). With the assembly sequence

**(a) Forces on Brick**   **(b) Structural Force Model**   **(c) Physics Reasoning**

Gravity.

Vertical normal and friction forces due to connections.

Horizontal normal forces due to connections.

Horizontal normal forces due to adjacent bricks.

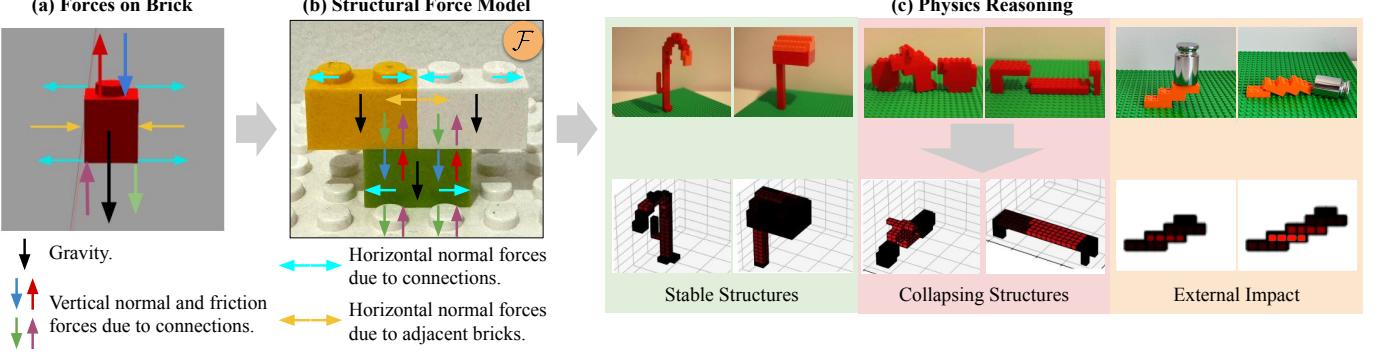Stable Structures   Collapsing Structures   External Impact

Fig. 3: **Assembly Physics Reasoning.** (a) All possible forces that could be exerted on a brick. (b) The force model $\mathcal{F}$ of an example structure. (c) Given a brick structure $B$ (top), the physics reasoning solves $S = f_{\text{stability}}(B)$ to determine each brick's stability score $s_i$ and evaluate the structural stability (bottom). Darker color indicates higher stability. Red: collapsing bricks, *i.e.*, $s_i = 0$.

$Q$, BRICKMATIC leverages APEX-MR [5], an asynchronous planning framework, to reason over 1) robot tasks (which robot performs which assembly step $\{b_i^a, a_i\}$), 2) robot motions (how to perform the assembly step without collision), and 3) collaborative executions (how to asynchronously collaborate to achieve better efficiency). Specifically, BRICKMATIC generates a bimanual robot execution plan $G$ from $\hat{B}$ as:

$$G = f_{\text{BRICKMATIC}}(\hat{B})$$
$$\text{s.t.} \quad B_i^a \in \mathcal{S}, \tag{3}$$
$$a_i \in \mathcal{A}, \ \forall i \in \{1, 2, \ldots, N\},$$

By executing $G$, BRICKMATIC turns the design $\hat{B}$ into a physical assembly product $B$. The details of BRICKMATIC are discussed in Sec. V.

### E. Prompt-to-Product Pipeline Summary

The entire Prompt-to-Product pipeline operates as follows:

$$\hat{B} = f_{\text{BRICKGPT}}(u) \in \mathcal{B}(u \mid \mathcal{I}, \mathcal{S})$$
$$G = f_{\text{BRICKMATIC}}(\hat{B})$$
$$B \leftarrow \mathbb{ROBOT\_EXECUTION}(G) \in \mathcal{B}(u \mid \mathcal{I}, \mathcal{A}, \mathcal{S})$$

This staged but interconnected architecture ensures that high-level intent is preserved while satisfying the constraints of physical realizability and robotic execution. The central role of the physics reasoning module in constraining design generation and guiding robot construction provides a critical coupling between semantic generation and physical grounding.

## III. ASSEMBLY PHYSICS REASONING

We leverage the stability analysis method in [13] to estimate the physical feasibility, *i.e.*, structural stability $S$, of any brick structure $B \in \hat{\mathcal{B}}$. Specifically, we optimize over force-balancing equations and solve for the force distribution to infer the stability of the brick structure.

**Structure Force Model.** For a brick structure $B$, we reason the possible forces that could be exerted on it (Fig. 3(a)), and construct the structural force model (Fig. 3(b)), which consists of a set of forces $\mathcal{F}$, depending on the connections. Each brick $b_i$ has $M_i$ forces exerting on it, where each force

$F_i^j \in \mathcal{F}_i \subseteq \mathcal{F}, j \in \{1, 2, \ldots, M_i\}$. A structure is stable if all bricks can reach static equilibrium:

$$\forall i \in \{1, 2, \ldots, N\},$$
$$\sum_j^{M_i} F_i^j = 0, \quad \sum_j^{M_i} \tau_i^j \doteq \sum_j^{M_i} L_i^j \times F_i^j = 0, \tag{4}$$

where $L_i^j$ denotes the force lever corresponding to $F_i^j$.

**Structural Stability Analysis.** To evaluate if a brick structure is physically stable, we solve its force distribution $\mathcal{F}$ by formulating the problem into a nonlinear program as

$$\underset{\mathcal{F}}{\operatorname{argmin}} \sum_i^N \left\{ |\sum_j^{M_i} F_i^j| + |\sum_j^{M_i} \tau_i^j| + \alpha \mathcal{D}_i^{\max} + \beta \sum \mathcal{D}_i \right\}, \tag{5}$$

subject to three constraints: 1) Non-negativity: $\forall F_i^j \in \mathcal{F}, F_i^j \geq 0$; 2) Non-coexistence: at any connecting point, the top brick cannot be supported while pulling by the bottom brick simultaneously; 3) Newton's third law: the forces exerted on adjacent bricks should be equal and opposite. $\mathcal{D}_i \subset \mathcal{F}_i$ is the set of friction forces that drag $b_i$ down (*i.e.*, the green arrow in Fig. 3(a)) and $\mathcal{D}_i^{\max}$ denotes the maximum value in $\mathcal{D}_i$. Hyperparameters $\alpha, \beta$ are taking the same values as in [13].

Solving the nonlinear program in Eq. (5) finds a force distribution $\mathcal{F}$ that drives the structure to static equilibrium with the minimum internal friction. From the solved $\mathcal{F}$, we obtain the structural stability $S$ and the score per brick as

$$s_i = \begin{cases} 0 & \begin{aligned} & \sum_j^{M_i} F_i^j \neq 0 \\ & \vee \sum_j^{M_i} \tau_i^j \neq 0 \\ & \vee \mathcal{D}_i^{\max} > F_T, \end{aligned} \\ \frac{F_T - \mathcal{D}_i^{\max}}{F_T} & \text{otherwise,} \end{cases} \tag{6}$$

where $F_T$ is a measured constant friction capacity between brick connections. Higher $s_i$ indicates greater stability, while $s_i = 0$ indicates an unstable brick that will cause structural failure: either the structure cannot reach static equilibrium (*i.e.*, $\sum_j^{M_i} F_i^j \neq 0 \vee \sum_j^{M_i} \tau_i^j \neq 0$) or the required friction exceeds the friction capacity of the material (*i.e.*, $\mathcal{D}_i^{\max} > F_T$). For a physically stable structure, we have $\forall i \in \{1, 2, \ldots, N\}, s_i > 0$. For any brick structure $B$, we can evaluate the structural
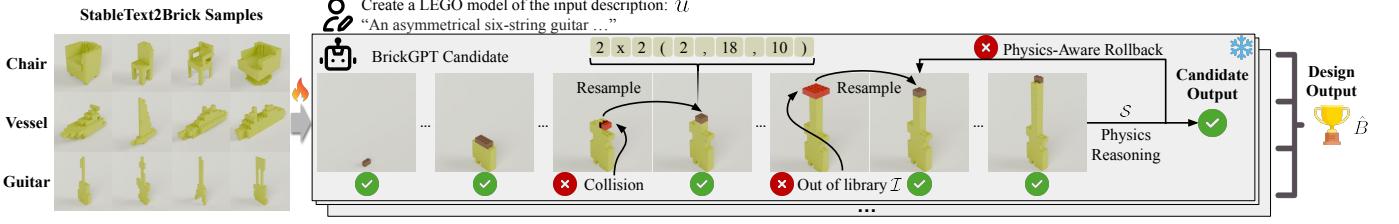
Fig. 4: **BRICKGPT for Generative Assembly Design.** BRICKGPT is fine-tuned on StableText2Brick [22]. At inference time, multiple BRICKGPT candidates generate multiple assembly designs in parallel, and the best design is output as the final design.

stability $S = f_{\text{stability}}(B)$ by solving Eq. (5) and (6).

Fig. 3(c) illustrates examples of stable and collapsing brick structures with their visualized stability scores (*i.e.*, static stability). In addition, the stability analysis can account for external impact (*i.e.*, dynamic stability), determining whether the brick structure can withstand external weight. The capability of numerically evaluating the structural stability mitigates the need for a reliable simulator. Consequently, Prompt-to-Product does not rely on a physics engine to understand real-world physics. Instead, Prompt-to-Product leverages 1) the static stability to constrain BRICKGPT to generate physically buildable brick assembly designs $\hat{B}$ and 2) the dynamic stability to guide BRICKMATIC to plan appropriate bimanual operations $G$ to construct the design.

## IV. BRICKGPT: GENERATIVE ASSEMBLY DESIGN

We leverage BRICKGPT [22] to generate physically buildable brick assembly structures from text prompts. BRICKGPT is an end-to-end method that skips intermediate representations and directly generates the assembly design from a user text prompt as indicated in Eq. (2). It is built based on a large language model (LLM) backbone to understand text prompt input. However, LLM generates text output instead of bricks. To ensure compatibility with LLM, BRICKGPT represents bricks in plain text as "$\{h\} \times \{w\}$ ($\{x\}, \{y\}, \{z\}$)", where $h, w$ are brick dimensions in the $X$, $Y$ directions, and $x, y, z$ are its coordinates. As shown in Fig. 4, a brick can be represented by 10 text tokens. Note that this representation is directly convertible to the representation in Sec. II, in which $h$ and $w$ collectively determine the brick type $c$ and orientation $\omega$, and $x, y, z$ are equivalent to $p$. Also, this representation encodes the brick in a more geometrically explicit way, which is easier for model learning and independent of the representation of $\mathcal{I}$. As a result, BRICKGPT generates a brick assembly by outputting text responses word by word. To generate physically buildable designs, BRICKGPT encodes physical constraints in both its training and inference phases.

### A. Encoding Physical Constraints in Model Fine-tuning

Despite having common knowledge, the LLM backbone lacks domain knowledge and the necessary geometric understanding for brick assembly. Thus, BRICKGPT is fine-tuned on a large-scale brick assembly dataset, *i.e.*, StableText2Brick [22], as shown in Fig. 4. Each structure is paired with different captions. Importantly, all brick structures in StableText2Brick are verified to be physically buildable using the physics reasoning in Eq. (1), *i.e.*, StableText2Brick $\subset \mathcal{B}(u \mid \mathcal{I}, \mathcal{S})$. BRICKGPT is fine-tuned using the stable text-brick pairs.

### B. Encoding Physical Constraints in Model Inference

Fine-tuning on a buildable dataset enables BRICKGPT to gain knowledge in generating brick designs. However, it remains difficult to ensure the generated assembly design is always compliant with the physical constraints. Thus, BRICKGPT further incorporates physical constraints into autoregressive inference. Specifically, BRICKGPT decouples the physical constraints and hierarchically enforces them.

**Brick-by-Brick Rejection Sampling.** During autoregressive inference, the constraints are relaxed to only consider resource and collision. As depicted in Fig. 4, after the model generates a brick $\hat{b}_i$, it immediately checks if $\hat{b}_i$ is valid, *i.e.*, 1) $\hat{c}_i \in \mathcal{I}$, and 2) $\hat{b}_i$ does not collide with the existing structure. If $\hat{b}_i$ is invalid, BRICKGPT rejects it and resamples a new brick.

**Physics-Aware Rollback.** Structural stability is applied at the end of the generation. Specifically, the model uses the physics reasoning in Eq. (1) to verify $\hat{B} \in \mathcal{S}$. Otherwise, BRICKGPT rolls back to a stable partial design, and regenerates from a partial design as shown in Fig. 4.
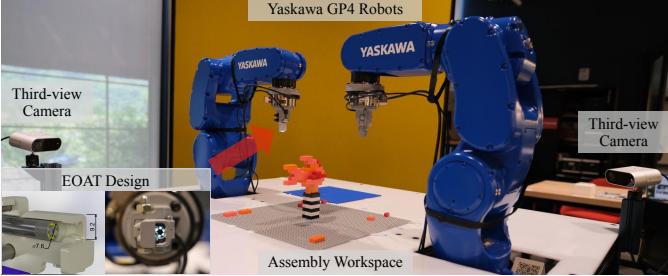
### C. Multi-Head Generation

Due to the probabilistic nature, BRICKGPT could generate different designs with an identical text prompt input. Thus, we present BRICKGPT++, which generates multiple designs in parallel to improve the design generation quality. Given a prompt input $u$, we run eight BRICKGPT instances with different seeds and generate multiple candidate designs. We render the outputs and select the best one according to the CLIP scores [23] as shown in Fig. 4.

By integrating physical constraints, *i.e.*, resource and stability, in both model fine-tuning and inference, BRICKGPT++ generates brick assembly designs $\hat{B}$ that satisfy the users' customization requirements and are physically buildable in reality, *i.e.*, $\hat{B} \in \mathcal{B}(u \mid \mathcal{I}, \mathcal{S})$.

## V. BRICKMATIC: BIMANUAL ASSEMBLY CONSTRUCTION

We present BRICKMATIC, a bimanual robotic system, to construct customized brick structures. In particular, this work presents a robotic system design that enables general-purpose robot arms to robustly and precisely manipulate individual bricks. In addition, a robot skill set $\mathcal{A}$ is introduced, which consists of a set of skill primitives, significantly expanding the system's dexterity. To construct a full brick structure, we further develop a multi-level reasoning framework, which generates an execution plan $G$ as defined in Eq. (3).

Fig. 5: **BRICKMATIC Embodiment Design.** Our BRICKMATIC includes two Yaskawa GP4 robots, each equipped with an FTS and an EiF EOAT, detailed in the bottom-left corner. A baseplate is placed in the middle of the workstation for the robots to construct brick assembly structures. Two third-view cameras are placed diagonally.

## A. Robot Embodiment Design

We present BRICKMATIC as shown in Fig. 5. BRICKMATIC is a bimanual robotic system, consisting of two Yaskawa GP4 robot arms. Each arm is equipped with an ATI force torque sensor (FTS) and a customized EOAT. Two cameras are mounted diagonally at the workstation. By facilitating dual arms, BRICKMATIC covers a larger workspace, achieves higher task efficiency, and most importantly, is capable of collaboratively performing more dexterous manipulation tasks.

BRICKMATIC leverages the Eye-in-Finger (EiF) [24] design as the EOAT to precisely and reliably manipulate bricks. The design of EiF is detailed in the bottom left of Fig. 5. The tooltip features a hollow interface that snaps over knobs on a brick and holds it from the top. It also has knobs on the side that allow it to securely hold a brick from the bottom. Most importantly, EiF integrates an endoscope camera inside the tooltip. This type of camera is selected due to its low cost, small form factor, and built-in LED, which ensures consistent lighting. By integrating the camera into the EOAT, EiF enables close-proximity visual feedback without increasing the geometric volume.

## B. Robot Skill Set

The multi-modality of BRICKMATIC enables the system to perform a wide variety of different operations. We parameterize the skill policies and construct a skill set $\mathcal{A}$ as shown in Fig. 6. Our skill set has three major categories of skills: 1) manipulation skills, 2) perception skills, and 3) motion skills.

**Manipulation Skills.** The manipulation skills represent structured, reusable action primitives that use sensor feedback (*e.g.*, force) to interact with the environment. Following the design in [5], our skill set includes six manipulation skills as shown in Fig. 6(a).

1) **Pick**: The robot disassembles a brick from the existing structure and picks it up. Due to the EiF design, we leverage the insert-and-twist policy [14], and the "Pick" skill is represented by two learnable parameters: a disassembly axis $O_d$ and a twisting angle $\theta_d$.
2) **Place-Down**: The robot assembles a brick initially in its hand (holding from the top) and places it at a target location. Similarly, we use the insert-and-twist policy, and the action is represented by two learnable parameters: an assembly axis $O_a$ and a twisting angle $\theta_a$.

3) **Place-Up**: The robot holds a brick from its bottom and assembles it in a location beneath another brick. We use the insert-and-twist policy, and the skill is represented by two learnable parameters: an assembly axis $O_a^u$ and a twisting angle $\theta_a^u$.
4) **Support-Bottom**: The robot supports the brick structure from the bottom.
5) **Support-Top**: The robot supports the brick structure from the top.
6) **Handover**: One robot is holding a brick from the top and places it onto the EiF EOAT of the other robot. After the operation, the other robot holds the brick from its bottom. This skill shares the same $O_a$ and $\theta_a$ from the "Place-Down" skill.

All manipulation skills are goal-conditioned, *i.e.*, conditioned on the goal location to support or pick/place a brick, and use FTS feedback to 1) avoid excessive pressure when picking and placing bricks, and 2) ensure a slight contact when supporting to prevent breaking the structure. We follow the learning technique in [14] to optimize the twisting axes $O_d, O_a, O_a^u$ and twisting angles $\theta_d, \theta_a, \theta_a^u$.

**Perception Skills.** BRICKMATIC can robustly manipulate individual bricks with the manipulation skills. However, contingencies could occur when building full structures. For instance, bricks assembled earlier could be loosened due to later operations, causing the structure to tilt, or even collapse. Thus, we introduce new perception skills, which leverage either the EiF in-hand camera or the third-view camera to update the BRICKMATIC's knowledge on the environment and detect if any contingency occurs while the robots remain in place. We present four perception skills as illustrated in Fig. 6(b).

1) **Detect Place**: The system detects if a place operation (*i.e.*, "Place-Down" or "Handover") is successful using the EiF camera. Essentially, this skill detects if the brick is released after the operation. Specifically, we learn a binary classifier on top of the DINOv2 features [19] of the EiF camera input. The classifier determines that the operation is successful if the brick is released. The "Detect Place" block in Fig. 6(b) depicts examples of camera views and their corresponding output.
2) **Detect Pick**: The system detects if a "Pick" skill is successfully performed using the EiF camera. The same classifier determines that the operation is successful if the brick is in-hand. The "Detect Pick" block in Fig. 6(b) depicts examples of detections.
3) **Detect Anomaly**: Due to the non-rigid connections between bricks, previously established connections could be loosened due to later operations, causing the structure to tilt or even collapse. To address this, we use third-view cameras as shown in Fig. 5 to detect anomaly, *i.e.*, structural failure. As shown in the "Detect Anomaly" block in Fig. 6(b), we compare the real camera view of the structure (*i.e.*, the bounding box) with the desired camera view rendered in simulation [9] (*i.e.*, the red visuals and their depth maps). If the real observation does not match the desired simulation rendering, then we classify the current scene as an anomaly.
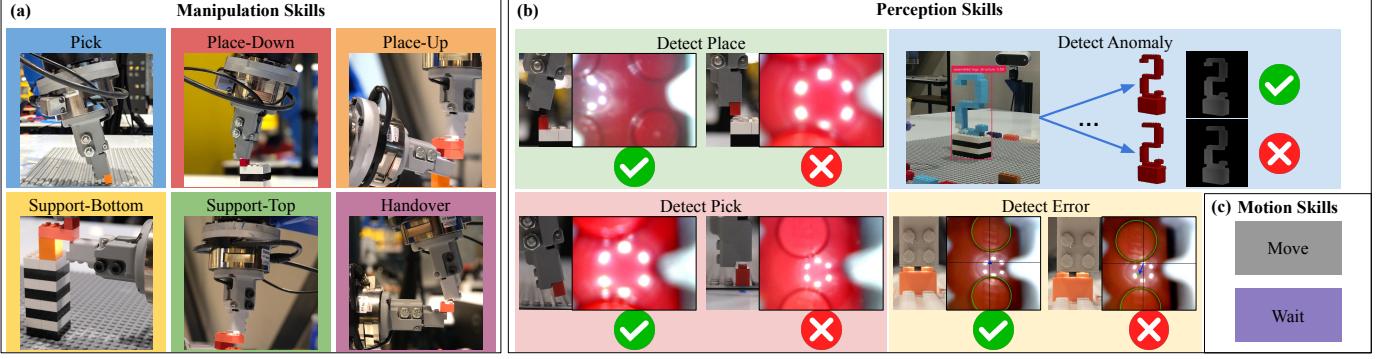
Fig. 6: **BRICKMATIC Skill Set.** The skill set $\mathcal{A}$ of BRICKMATIC with the presented embodiment.

4) **Detect Error**: The system's knowledge of the environment could deviate from reality due to, for instance, calibration error, tilted structures, etc. Following [24], we use the endoscopic camera to detect and correct model mismatches. As shown in the "Detect Error" block in Fig. 6(b), we use a fine-tuned YOLOv8-seg [7] to detect brick knobs and calculate the offset between the actual brick position and the center of EiF. The robot adjusts its pose accordingly to ensure accurate alignment with the brick, which is essential for subsequent manipulation skills, *e.g.*, "Pick", "Place-Down", etc.

**Motion Skills.** The motion skills are skills that do not physically interact with the environment. Our skill set includes two fundamental motion skills: "Wait" and "Move" shown in Fig. 6(c). The "Move" skill is responsible for collision-free movements between waypoints using motion planners, *e.g.*, [10, 12]. The "Wait" skill pauses the robot's motion.

### C. Multi-Level Reasoning

The embodiment design and the skill set provide BRICK-MATIC enhanced dexterity to reliably manipulate individual bricks. To construct a complete brick structure, we present a multi-level reasoning framework to plan and coordinate the bimanual system. As illustrated in Fig. 7, the presented multi-level reasoning framework considers the embodiment design, the skill set $\mathcal{A}$, and the brick design $\hat{B}$, and generates an execution plan $G$ that controls the bimanual system to construct the physical brick assembly structure $B \in \mathcal{B}(u \mid \mathcal{I}, \mathcal{A}, \mathcal{S})$.

**Assembly Reasoning.** Even though BRICKGPT generates the brick design brick-by-brick, it only ensures the final design $\hat{B}$ is physically buildable. Hence, intermediate structures could be unstable as illustrated in Fig. 2(a). Thus, given the brick assembly design $\hat{B}$, the first step is to plan a physically executable assembly sequence by reordering the bricks as depicted in Fig. 2(b). Specifically, we use assembly-by-disassembly search [25] to search for an assembly sequence. As illustrated in Fig. 7(a), we start from the goal state, *i.e.*, the final assembly design, and disassemble a brick step by step until there is no brick remaining. The assembly sequence is obtained by reversing the disassembly sequence.

However, for each disassembly step, multiple bricks could be disassembled, and it is crucial to choose the physically executable one so that the planned assembly sequence is executable by BRICKMATIC. To this end, we leverage the action mask presented in [15] to evaluate if a brick $\hat{b}_i$ is removable from a partial structure $\hat{B}_i$ using skill $a_i \in \mathcal{A}$ by assessing the following properties:

1) **Operability**: there is space for the system to perform the skill $a_i$. Let $\mathcal{V}_{a_i}$ be the volume occupied by $a_i$, we ensure $\mathcal{V}_i \cap \mathcal{V}_{a_i} = \varnothing$, where $\mathcal{V}_i$ denotes the volume occupied by $\hat{B}_i$. Since $a_i$ is a goal-conditioned skill primitive, we can easily estimate the occupied space.
2) **Static stability**: the structure after removing the brick remains physically stable. Let $\hat{B}_i'$ be the structure resulted from removing $\hat{b}_i$ from $\hat{B}_i$, we ensure $\hat{B}_i' \in \mathcal{S}$.
3) **Dynamic stability**: the structure $\hat{B}_i$ is stable under the impact of $a_i$. Since the physics reasoning $f_{\text{stability}}(\cdot)$ can account for external impact by adding virtual bricks as shown in Fig. 3, we model the robot place operation as a virtual brick $\hat{b}_p$ with heavy mass, *i.e.*, 1kg. If support is needed, we add an additional virtual brick $\hat{b}_s$ with mass being $-1$kg. We apply the physics reasoning to the virtual structure $\hat{B}_i^+ = \hat{B}_i \cup \{\hat{b}_p, \hat{b}_s\}$ and ensure $\hat{B}_i^+ \in \mathcal{S}$.

By ensuring the above criteria, we ensure the disassembly step $\{\hat{b}_i, a_i\}$ is physically executable. Note that $a_i$ can be a place skill alone or in combination with a support skill. We use the action mask to prune any invalid disassembly steps, which are shown as dashed arrows in Fig. 7(a). Among these criteria, static and dynamic stability evaluations with physics reasoning in Eq. (1) are more expensive than evaluating operability. Thus, we adopt a modified Depth-First Search (DFS) with partial expansion and parallel action mask evaluation. At each node (*i.e.*, a partial structure), we identify all operable bricks and only sample a subset for parallel stability evaluation. If no executable actions are found, we resample until all options are exhausted, then backtrack. This partial expansion approach enables efficient discovery of a physically executable assembly sequence $Q$.

Given an assembly sequence $Q$, we generate a robot executable plan $G$ following the technique presented in APEX-MR [5] as shown in Fig. 7(b)-(d).

**Task Reasoning.** Following the idea in APEX-MR, we start from a sequential task plan, in which only one robot is moving at a time, as depicted in Fig. 7(b). We distribute the tasks $\{b_i^a, a_i\}$ to different robots by solving an integer-linear program and construct task dependencies, *i.e.*, the gray dashed arrow, indicating that a certain skill (pointing to) should be performed
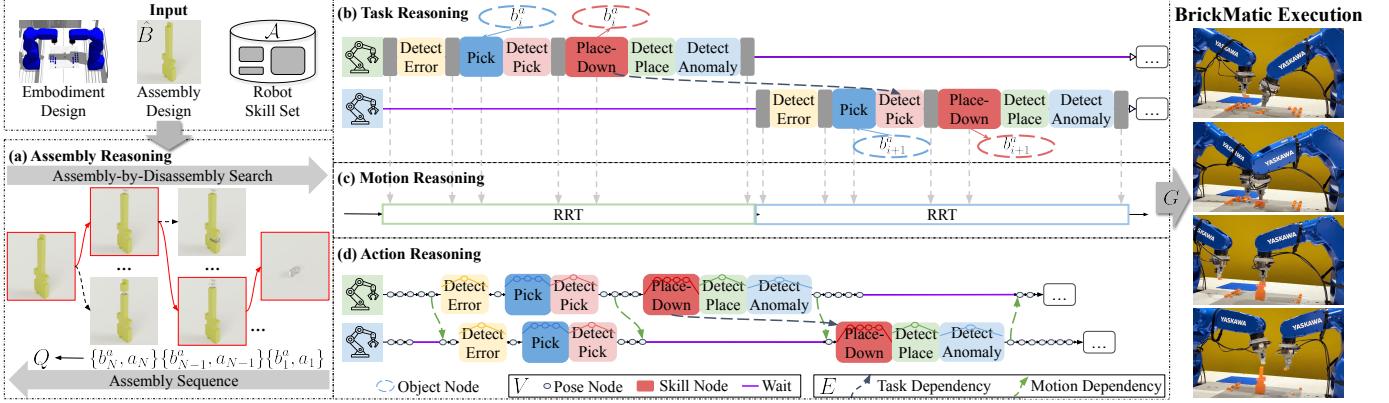
Fig. 7: **BRICKMATIC Multi-Level Reasoning for Bimanual Assembly Construction.** Given an assembly design, BRICKMATIC (a) reasons the assembly structure and generates a physically executable step-by-step assembly task plan $Q$. It then (b) reasons the tasks, distributes the assembly steps to different robots, and builds a sequential task plan using the available skills from the skill set $\mathcal{A}$. After that, BRICKMATIC does (c) motion reasoning to plan the motions for each robot to perform the tasks, and (d) action reasoning to convert the sequential motion plan into a parallel, asynchronous action plan, *i.e.*, a TPG $G = \{V, E\}$, for physical execution. The colors of the skill nodes correspond to the colors in the skill set in Fig. 6.

after the other (pointing from). In addition to manipulation skills discussed in APEX-MR, we add perception skills before and after manipulation skills to improve their robustness and detect failures.

**Motion Reasoning.** Given the sequential task plan, the robot motions (*i.e.*, the pose nodes) for each robot to perform the skills are planned using RRT-Connect [10]. Since only one robot is moving at a time, the bimanual cooperation is simplified to a single-agent planning problem, which significantly reduces the planning time and improves the quality of the generated collision-free motion plan.

**Action Reasoning.** Given the collision-free motion plan, we iterate over the pose nodes and construct motion dependencies, *i.e.*, the green dashed arrows, by identifying colliding node pairs. In the end, we construct a TPG with the nodes $V$ and dependency edges $E$, *i.e.*, $G = \{V, E\}$, which allows the robots to execute in parallel and collaborate asynchronously. Each robot only waits when necessary, significantly improving collaboration efficiency while avoiding potential collisions. Note that, unlike APEX-MR, which only considers force feedback, we incorporate perception skills before and after manipulation skills to improve their robustness and detect failures. When a failure is detected, the perception skill blocks the TPG execution and keeps detecting until the failure is recovered, *e.g.*, addressed by a human operator. The system continues, and the TPG execution flow remains intact since all precedence constraints are still satisfied.

By integrating physical constraints, *i.e.*, dexterity and feasibility, BRICKMATIC generates bimanual robot operations $G$. Executing $G$ enables the system to efficiently, safely, and robustly construct customized brick assembly structures, bringing virtual brick designs $\hat{B}$ to real products $B \in \mathcal{B}(u \mid \mathcal{I}, \mathcal{A}, \mathcal{S})$.

## VI. EXPERIMENT

### A. Prompt-to-Product Interface

The Prompt-to-Product is hosted on a server with eight Nvidia RTX A4000 GPUs. We create a web interface to allow

TABLE I: **Generative Design Comparison.** Results are averaged over 36 open-world unique prompts from the user study in Sec. VI-C. % Buildable: the percentage of prompts with a generated BRICKMATIC-buildable brick design. CLIP: text-image similarity.

| Method | % Buildable | CLIP [23] | Time (s) |
|---|---|---|---|
| BRICKGPT [22] | 19.4% | 0.248±0.036 | **44.0±39.0** |
| BRICKGPT++ | **66.6%** | **0.266±0.029** | 85.7±48.8 |

users to interact with Prompt-to-Product as shown in Fig. 9(d). On the front page, users enter the text prompt in the prompt window and BRICKGPT++ generates the assembly designs. Due to the multi-head generation, the system outputs multiple designs and highlights the best one. Users then select one or more designs based on their preference and proceed with BRICKMATIC. In our interface, users do not directly interact with the physical BRICKMATIC system. Instead, we develop a BRICKMATIC digital twin in Gazebo [9], and the web interface displays BRICKMATIC running virtually. Note that since the digital twin cannot simulate the connections between bricks, we turn off physical interactions between bricks. Thus, the virtual BRICKMATIC only uses the manipulation and motion skills, and the execution is deterministic. After users verify their preferred designs virtually, we run BRICKMATIC offline, which fully utilizes the skill set $\mathcal{A}$, to physically construct the brick assembly products. Due to the staged architecture of Prompt-to-Product, users can also interact with BRICKGPT or BRICKMATIC individually. They can query BRICKGPT to generate a design and manually build it as depicted in Fig. 9(b), or they can input a manual brick design and let BRICKMATIC automatically construct it as shown in Fig. 9(c).

### B. Module Evaluation

**Generative Assembly Design.** The multi-head generation in Sec. IV-C improves the quality of design generation. Table I shows the comparison between the modified BRICKGPT++ and the original BRICKGPT [22], which generates one design per prompt. As shown by the buildable rate, the modified

TABLE II: **Bimanual Construction Comparison.** Success Rate: number of trials attempted to have the system successfully build the brick design once without restarting. Survival Length: number of bricks (averaged over the attempts) the system assembled without restarting. Time: the time to plan the assembly sequence $Q$, and construct TPG $G$ given the brick design.

| Method | Design | Success Rate | Survival Length | Time (s) |
|---|---|---|---|---|
| Dual-Arm [5] | Faucet | 1/5 | 9.2 | 30.0 |
| | Fish | 0/5 | 7.8 | 157.0 |
| | Vessel | 1/3 | 33.7 | 78.4 |
| | Guitar | **1/1** | **24** | 50.8 |
| BRICKMATIC | Faucet | **1/1** | **14** | **27.2** |
| | Fish | **1/1** | **29** | **83.4** |
| | Vessel | **1/1** | **36** | **58.1** |
| | Guitar | **1/1** | **24** | **42.9** |



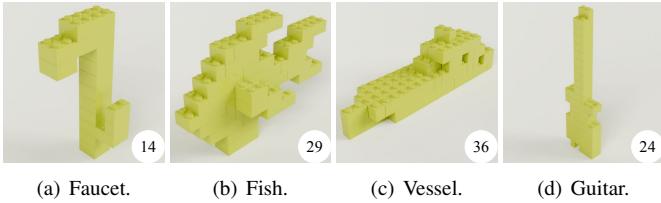| (a) Faucet. | (b) Fish. | (c) Vessel. | (d) Guitar. |

Fig. 8: Brick assembly designs for comparing bimanual assembly construction. The number in each figure indicates the number of bricks required.

BRICKGPT++ has a higher chance of generating BRICKMATIC-buildable designs. In addition, the CLIP score indicates that the designs from BRICKGPT++ better align with the open-world user prompts. Despite the longer generation time, BRICKGPT++ improves the quality of the design generation and enhances the overall Prompt-to-Product workflow.

**Bimanual Assembly Construction** We compare the presented BRICKMATIC with the dual-arm system in [5]. Specifically, BRICKMATIC is equipped with innovative EiF EOATs providing close-proximity visual feedback and has an expanded skill set $\mathcal{A}$ with perception skills for anomaly detection. Fig. 8 illustrates brick assembly designs we use to compare the systems. Here, we define a successful build as a trial in which the system builds the structure without manually stopping or restarting. As shown in Table II, with the integration of perception skills, BRICKMATIC has a significantly higher success rate as it successfully builds all designs in one attempt. This is because when a failure happens, BRICKMATIC automatically pauses its construction and actively requests human intervention for failure recovery. It continues its construction when the failure is addressed without starting over. On the other hand, the dual-arm system [5] needs to restart the entire building sequence once a failure happens. Consequently, it requires significantly more attempts and even fails to successfully build the Fish after many trials. In addition, the survival length reveals that BRICKMATIC can perform significantly longer horizon assembly tasks without restarting. Moreover, we compare the planning time for both systems. Due to the modified DFS in Sec. V-C, BRICKMATIC achieves a significantly shorter runtime.

### C. User Study

We conduct an IRB-approved user study to investigate how the proposed Prompt-to-Product can help people create physical assemblies from abstract ideas. Specifically, we want to test the following hypothesis:

1) Generative AI (BRICKGPT) helps users to bring their abstract ideas into concrete assembly designs. It reduces the required manual effort and expert knowledge.
2) The designs from BRICKGPT are good initial designs for users to improve further.
3) Bimanual robotic assembly (BRICKMATIC) reduces the required manual effort in creating physical products from virtual designs.
4) Users would prefer to build by themselves for fun if they were only to build a few assemblies. They would prefer robotic construction for mass production.

**Tasks.** Participants were divided into three groups: BRICK-GPT-only, BRICKMATIC-only, and the full Prompt-to-Product. In the BRICKGPT-only group, users submitted a prompt, created a manual design, and compared it with the best generated result (Fig. 9(b)). In the BRICKMATIC-only group, users assembled a given design with 3D visualization support (Fig. 9(c)), then watched a video of the BRICKMATIC performing the same task and rated their preference. In the full Prompt-to-Product group, users entered a prompt, created a manual design, selected a generated design, and compared the manual design with the one assembled in simulation by BRICKMATIC (Fig. 9(d)).

**Participants.** We recruited 21 participants (ages 18–31; 6 female and 15 male), primarily with technical backgrounds, while some came from the arts or architecture fields. LEGO experience ranged from none (2 participants) to advanced (4 had built sets >1000 pieces). In our study, 8 users tried BRICKGPT, 7 used BRICKMATIC, and 6 interacted with the full pipeline.

**Procedure.** Each user interacted with our system for ∼30 minutes. For BRICKGPT-only, each user had 4 design attempts, including 2 provided prompts and 2 customized prompts in one of ten categories (Bookshelf, Car, Chair, Guitar, Sofa, Table, Vessel, Bench, Bottle, Bus). For BRICKMATIC-only, each user assembled 2-4 brick designs, depending on the time needed to replicate the exact design. For the entire Prompt-to-Product pipeline, each user tried two customized prompts from two categories (Guitar, Vessel). All participants filled out a post-study survey, rating the system usability and favorability.

**Results.** User responses to key survey questions are summarized in Fig. 10. Overall, the Prompt-to-Product system is well-received and notably reduces user effort, especially when assembling multiple designs. For Hypotheses 1 and 3, Wilcoxon signed-rank tests on 5-point Likert responses show significant reductions in both physical and mental effort when using Prompt-to-Product compared to manual creation. BRICKGPT significantly reduces physical (p = 0.037) and mental effort (p = 0.041). BRICKMATIC yields a significant reduction in physical effort (p = 0.023) and a marginal reduction in mental effort (p = 0.055). The full pipeline shows the strongest effects, with highly significant reductions in both physical and mental
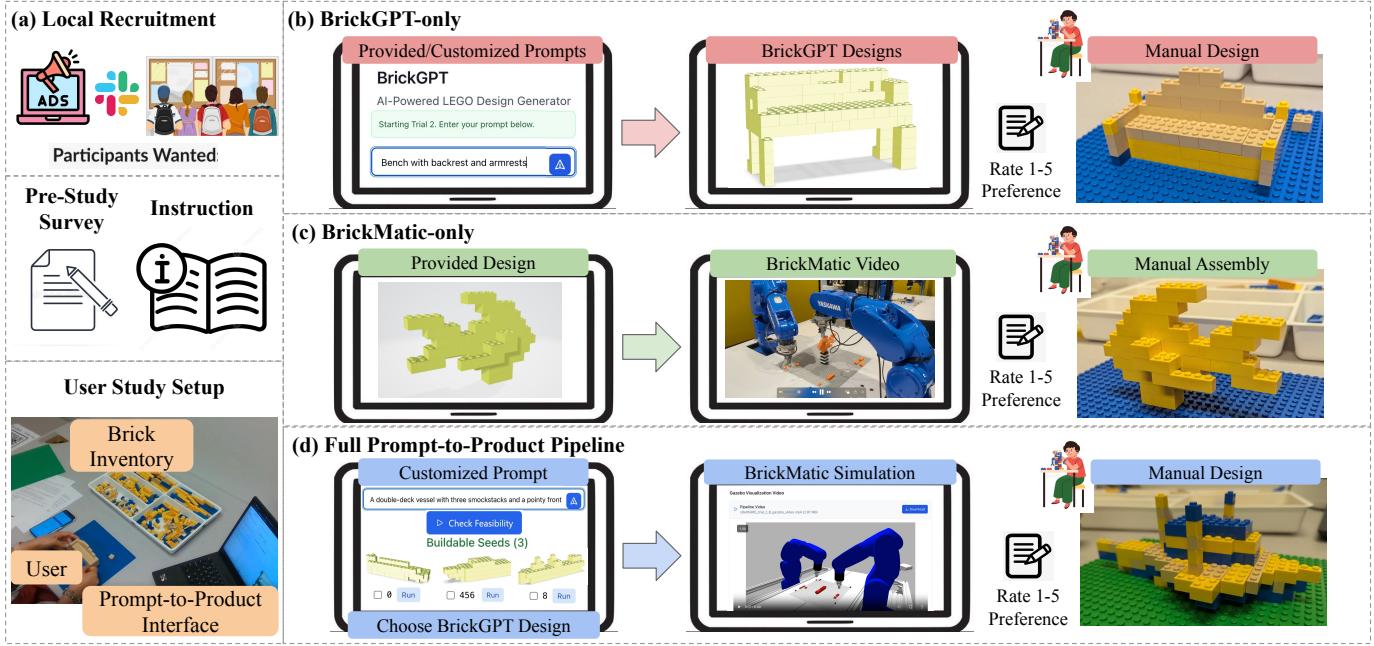
Fig. 9: **Prompt-to-Product User Study Procedure**. We recruit participants through public online advertisements, internal Slack channels, and physical bulletin boards. Each user completes a pre-study survey and interacts with our system. They also fill out a post-study survey in the end.

effort (p = 0.021 for both). For Hypothesis 4, users prefer using BRICKMATIC (p = 0.008) or the full pipeline (p = 0.021) when building multiple structures, but slightly prefer manual assembly for a single design with BRICKMATIC (p = 0.109) or the full pipeline (p = 0.039). For Hypothesis 2, no significant preference is observed between manual and BRICKGPT-generated designs as starting points (p = 0.5). However, with the full pipeline, users show a slight preference for generated designs (p = 0.097), suggesting they are at least equally suitable—or potentially better—for refinement.

Manual creation time ranged from under 2 to over 7 minutes, reflecting differences in user experience with brick assembly. Users generally preferred their manual designs over generated ones for custom prompts (3.79 vs. 2.86), while slightly preferred generated designs for provided prompts (3.44 vs. 3.50). Custom prompts were often highly creative (*e.g.*, a five-wheeled car, a tree-shaped bookshelf), challenging the system. Notably, all users expressed strong interest in using the Prompt-to-Product system again.

## VII. DISCUSSION

We propose Prompt-to-Product, a pipeline that translates abstract text prompts into physical brick assemblies, reducing required manual effort and expert knowledge. While promising, Prompt-to-Product has limitations to be address in future work.

**Generative Design:** Prompt-to-Product currently supports only brick assemblies. While users can modify the brick inventory $\mathcal{I}$, non-brick components are not supported, limiting the system from generating more vivid designs. Additionally, it is restricted to categories in StableText2Brick, leading to degraded performance on open-ended prompts—a limitation frequently noted by users during the study. Hence, we will explore more diverse 3D datasets and expand the generative

capabilities to include varied components and higher-fidelity assemblies that better match user intent.
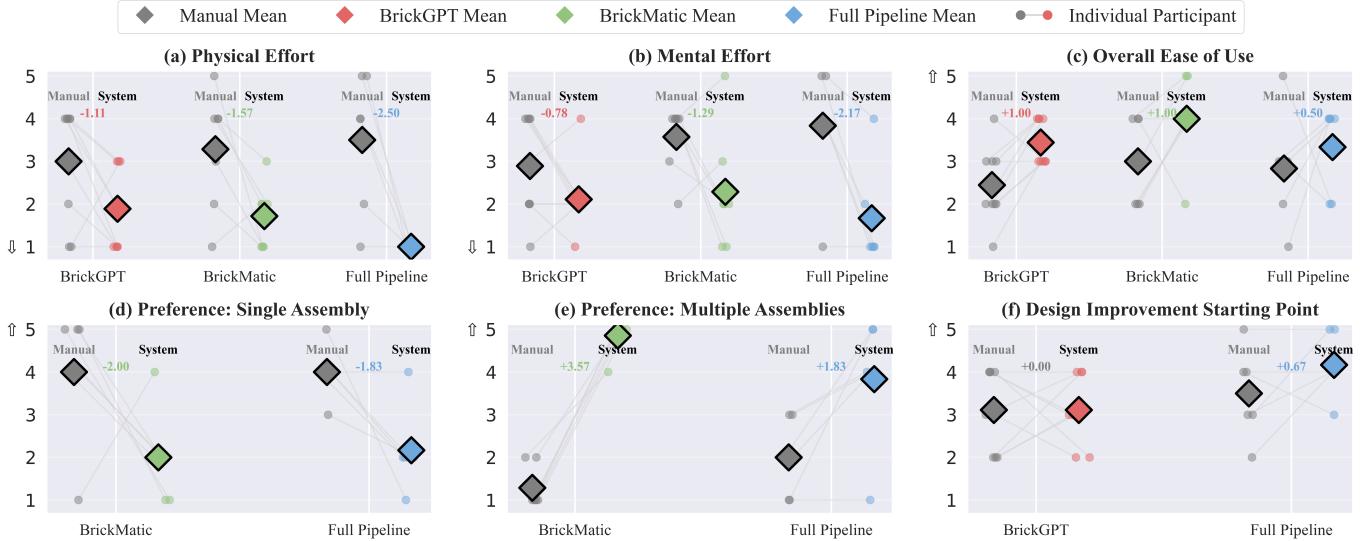
**System Dexterity:** Prompt-to-Product showcases state-of-the-art capabilities in long-horizon and dexterous manipulation tasks, but still falls short when compared to human-level dexterity, *i.e.*, user-constructable structures remain unbuildable by the current system (Table I). Thus, we will enhance BRICKMATIC by incorporating additional skills, such as in-hand assembly, failure recovery, and coordinated subassembly.

## REFERENCES

[1] Bambu Lab. Carbon X1, 2022. URL https://bambulab.com/en/x1.

[2] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. $\pi_0$: A vision-language-action flow model for general robot control. *arXiv:2410.24164*, 2024.

[3] Andrew Goldberg, Kavish Kondap, Tianshuang Qiu, Zehan Ma, Letian Fu, Justin Kerr, Huang Huang, Kaiyuan Chen, Kuan Fang, and Ken Goldberg. Blox-net: Generative design-for-robot-assembly using vlm supervision, physics, simulation, and a robot with reset. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2025.

[4] Wolfgang Hoenig, T. K. Kumar, Liron Cohen, Hang Ma, Hong Xu, Nora Ayanian, and Sven Koenig. Multi-agent path finding with kinematic constraints. In *International Conference on Automated Planning and Scheduling (ICAPS)*, 2016.

[5] Philip Huang, Ruixuan Liu, Shobhit Aggarwal, Changliu Liu, and Jiaoyang Li. Apex-mr: Multi-robot asynchronous planning

Fig. 10: **Prompt-to-Product User Study Results.** Individual and mean ratings of key survey questions for BrickGPT-only, BrickMatic-only, and entire Prompt-to-Product. (a) Physical demands of designing or building brick assembly structures. (b) Mental demands of designing or building brick assembly structures. (c) The overall ease of use. (d) User's preference when only building one brick design. (e) User's preference when building multiple designs. (f) User's rating of how well the manual or BRICKGPT-generated design serves as a starting point for improvement.

and execution for cooperative assembly. In *Robotics: Science and Systems (RSS)*, 2025.

[6] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, et al. $\pi_{0.5}$ : a vision-language-action model with open-world generalization. *arXiv:2504.16054*, 2025.

[7] G. Jocher, A. Chaurasia, and J. Qiu. Ultralytics YOLOv8, 2023. URL https://github.com/ultralytics/ultralytics.

[8] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv:2406.09246*, 2024.

[9] N. Koenig and A. Howard. Design and use paradigms for gazebo, an open-source multi-robot simulator. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2004.

[10] J.J. Kuffner and S.M. LaValle. Rrt-connect: An efficient approach to single-query path planning. *IEEE International Conference on Robotics and Automation*, 2000.

[11] Alexander Htet Kyaw, Se Hwan Jeon, Miana Smith, and Neil Gershenfeld. Speech to reality: On-demand production using natural language, 3d generative ai, and discrete robotic assembly. *arXiv:2409.18390*, 2024.

[12] Ruixuan Liu, Rui Chen, Yifan Sun, Yu Zhao, and Changliu Liu. Jerk-bounded position controller with real-time task modification for interactive industrial robots. In *2022 IEEE/ASME International Conference on Advanced Intelligent Mechatronics (AIM)*, pages 1771–1778, 2022. doi: 10.1109/AIM52237.2022.9863251.

[13] Ruixuan Liu, Kangle Deng, Ziwei Wang, and Changliu Liu. Stablelego: Stability analysis of block stacking assembly. *IEEE Robotics and Automation Letters (RA-L)*, 9(11):9383–9390, 2024.

[14] Ruixuan Liu, Yifan Sun, and Changliu Liu. A lightweight and transferable design for robust lego manipulation. *International Symposium on Flexible Automation (ISFA)*, 2024.

[15] Ruixuan Liu, Alan Chen, Weiye Zhao, and Changliu Liu. Physics-aware combinatorial assembly sequence planning using data-free action masking. *IEEE Robotics and Automation Letters (RA-L)*, 10(5):4882–4889, 2025.

[16] Sheng-Jie Luo, Yonghao Yue, Chun-Kai Huang, Yu-Huan Chung, Sei Imai, Tomoyuki Nishita, and Bing-Yu Chen. Legolization: optimizing lego designs. *ACM Transactions on Graphics (TOG)*,

November 2015.

[17] Vihaan Misra, Peter Schaldenbrand, and Jean Oh. Shapeshift: Towards text-to-shape arrangement synthesis with content-aware geometric constraints. *arXiv:2503.14720*, 2025.

[18] Ludwig Nägele, Alwin Hoffmann, Andreas Schierl, and Wolfgang Reif. Legobot: Automated planning for coordinated multi-robot assembly of lego structures. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020.

[19] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv:2304.07193*, 2023.

[20] Martin Pletz and Matthias Drvoderic. Brickfem an automated finite element model for static and dynamic simulations of simple lego® sets.

[21] Ivaylo Popov, Nicolas Heess, Timothy Lillicrap, Roland Hafner, Gabriel Barth-Maron, Matej Vecerik, Thomas Lampe, Yuval Tassa, Tom Erez, and Martin Riedmiller. Data-efficient deep reinforcement learning for dexterous manipulation. *arXiv*, 2017.

[22] Ava Pun, Kangle Deng, Ruixuan Liu, Deva Ramanan, Changliu Liu, and Jun-Yan Zhu. Generating physically stable and buildable brick structures from text. In *IEEE International Conference on Computer Vision (ICCV)*, 2025.

[23] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021.

[24] Zhenran Tang, Ruixuan Liu, and Changliu Liu. Eye-in-finger: Smart fingers for delicate assembly and disassembly of lego. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2025.

[25] Yunsheng Tian, Jie Xu, Yichen Li, Jieliang Luo, Shinjiro Sueda, Hui Li, Karl D.D. Willis, and Wojciech Matusik. Assemble them all: Physics-based planning for generalizable assembly by disassembly. *ACM Transactions on Graphics (TOG)*, 2022.

[26] Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, , Adam Letts, Yangguang Li, Ding Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. Triposr: Fast 3d object reconstruction from a single image. *arXiv:2403.02151*, 2024.

[27] Pinar Urhal, Andrew Weightman, Carl Diver, and Paulo Bartolo. Robot assisted additive manufacturing: A review. *Robotics and Computer-Integrated Manufacturing*, 59:335–345, 2019. ISSN 0736-5845.

[28] Longwen Zhang, Qixuan Zhang, Haoran Jiang, Yinuo Bai, Wei Yang, Lan Xu, and Jingyi Yu. Bang: Dividing 3d assets via generative exploded dynamics. *ACM Transactions on Graphics (TOG)*, 2025.

[29] Tianqi Zhang, Zheng Wu, Yuxin Chen, Yixiao Wang, Boyuan Liang, Scott Moura, Masayoshi Tomizuka, Mingyu Ding, and Wei Zhan. Physics-aware robotic palletization with online masking inference. *arXiv:2502.13443*, 2025.

[30] Zibo Zhao, Zeqiang Lai, Qingxiang Lin, Yunfei Zhao, Haolin Liu, Shuhui Yang, Yifei Feng, Mingxin Yang, Sheng Zhang, Xianghui Yang, Huiwen Shi, Sicong Liu, Junta Wu, Yihang Lian, Fan Yang, Ruining Tang, Zebin He, Xinzhou Wang, Jian Liu, Xuhui Zuo, Zhuo Chen, Biwen Lei, Haohan Weng, Jing Xu, Yiling Zhu, Xinhai Liu, Lixin Xu, Changrong Hu, Shaoxiong Yang, Song Zhang, Yang Liu, Tianyu Huang, Lifu Wang, Jihong Zhang, Meng Chen, Liang Dong, Yiwen Jia, Yulin Cai, Jiaao Yu, Yixuan Tang, Hao Zhang, Zheng Ye, Peng He, Runzhou Wu, Chao Zhang, Yonghao Tan, Jie Xiao, Yangyu Tao, Jianchen Zhu, Jinbao Xue, Kai Liu, Chongqing Zhao, Xinming Wu, Zhichao Hu, Lei Qin, Jianbing Peng, Zhan Li, Minghui Chen, Xipeng Zhang, Lin Niu, Paige Wang, Yingkai Wang, Haozhao Kuang, Zhongyi Fan, Xu Zheng, Weihao Zhuang, YingPing He, Tian Liu, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, Jingwei Huang, and Chunchao Guo. Hunyuan3d 2.0: Scaling diffusion models for high resolution textured 3d assets generation, 2025. URL https://arxiv.org/abs/2501.12202.