

Generative Factor Chaining: Coordinated Manipulation with Diffusion-based Factor Graph

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Learning to plan for multi-step, multi-manipulator tasks is notoriously
2 difficult because of the large search space and the complex constraint satisfaction
3 problems. We present Generative Factor Chaining (GFC), a composable genera-
4 tive model for planning. GFC represents a planning problem as a spatial-temporal
5 factor graph, where nodes represent objects and robots in the scene, spatial factors
6 capture the distributions of valid relationships among nodes, and temporal factors
7 represent the distributions of skill transitions. Each factor is implemented as a
8 modular diffusion model, which are composed during inference to generate feasi-
9 ble long-horizon plans through bi-directional message passing. We show that GFC
10 can solve complex bimanual manipulation tasks and exhibits strong generalization
11 to unseen planning tasks with novel combinations of objects and constraints. More
12 details can be found at: sites.google.com/view/generative-factor-chaining

13 **Keywords:** Manipulation Planning, Bimanual Manipulation, Generative Models

14 1 Introduction

15 Solving real-world sequential manipulation tasks requires reasoning about sequential dependencies
16 among manipulation steps. For example, a robot needs to grip the center or the tail of a hammer,
17 instead of its head, in order to subsequently hammer a nail. The complexity of planning problems
18 increases when multiple manipulators are involved, where spatial coordination constraints among
19 manipulators need to be satisfied. In the example shown in [Figure 1](#), the robot has to reason about
20 the optimal pose to grasp the hammer with the left arm, such that the right arm can coordinate to
21 re-grasp. Subsequently, the two arms must coordinate to hammer the nail. While classical Task
22 and Motion Planning (TAMP) methods have shown to be effective at solving such problems by
23 hierarchical decomposition [1], they require accurate system state and kinodynamic model. Further,
24 searching in such a large solution space to satisfy numerous constraints poses a severe scalability
25 challenge. In this work, we aim to develop a learning-based planning framework to tackle complex
26 manipulation tasks with both sequential and spatial coordination constraints.

27 To solve complex sequential manipulation problems, prior learning-to-plan methods have largely
28 adopted the options framework and modeled the preconditions and effect of the options or primitive
29 skills [2, 3, 4, 5, 6, 7]. Key to their successes are skill chaining functions that determine whether
30 executing a skill can satisfy the precondition of the next skill in the plan, and eventually the success
31 condition of the overall task. However, the use of vectorized states and the assumption of a linear
32 chain of sequential dependencies limits the expressiveness of these methods. Consider a task where
33 a robot fetches two items from a box. Intuitively, the skills for fetching one object should not
34 influence the other. However, due to vectorized states and the linear dependency assumption, the
35 skill-chaining methods are forced to model such sequential dependencies. Similarly, a skill intended
36 to satisfy a future skill’s condition will be forced to influence the steps in between. Finally, the skill
37 chain representation forbids these methods from effectively modeling multiple-arm manipulation
38 tasks, where concurrent skills must be planned to jointly satisfy a constraint.

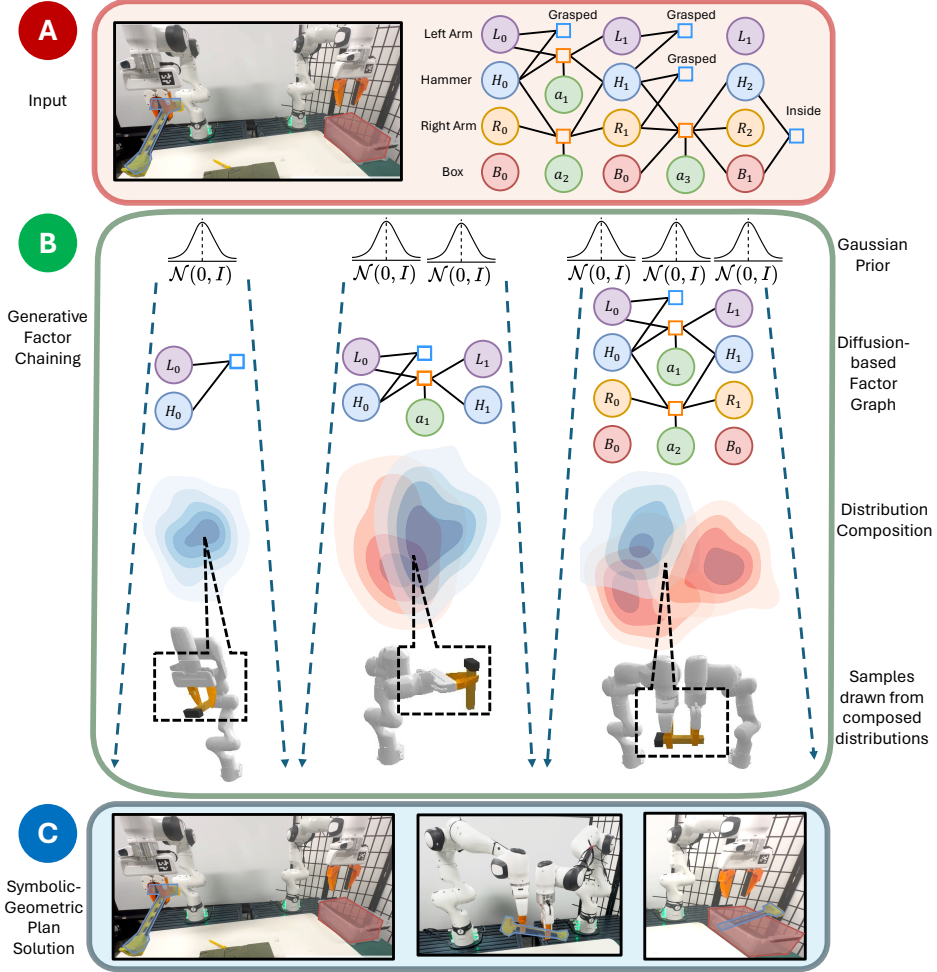


Figure 1: **Factor graph for a multi-arm coordination task.** Our factor graph-based planning formulation solves for a sequence of spatial factor graphs from the initial state to a goal factor by chaining them using temporal skill factors. The figure illustrates the temporal evolution of a factor graph by executing single or multiple skills sequentially or in-parallel to handover a hammer, pick up a nail, and coordinate both arms to strike the nail. **Task:** The task objective is to place the hammer inside the box. However, since the left arm cannot reach the box, the hammer is handed over to the right arm such that the right arm can complete the task. **(a) Inputs:** The initial scene and a symbolically feasible spatial-temporal factor graph plan to complete the goal objective. **(b) GFC:** We formulate all factors as distributions of the nodes connected to them. GFC represents spatial factors as classifiers and temporal factors as diffusion models. We leverage compositionality of diffusion models to compose spatial-temporal distributions and find the joint distribution of the complete plan directly at inference. Finally, samples drawn from such a joint distribution are symbolically and geometrically feasible solutions of the whole plan. **(c) Output:** A sequence of skill choices and optimizer continuous parameters executed on robots with parameterized skill controllers.

39 To move beyond the linear chain and model complex coordinated manipulation, we introduce Generative
40 Factor Chaining (GFC), a learning-to-plan framework built on flexible composable generative
41 models. For a given symbolically feasible plan graph, GFC adopts a spatial-temporal factor graph [8]
42 representation, where nodes are objects and robot states, and spatial factors represent the relationship
43 constraints between these nodes. Skills are temporal factors that connect these state-factor graphs
44 via transition distributions. A single skill factor can simultaneously connect to multiple object and
45 robot nodes, allowing for natural representation of complex multi-object interactions and steps that
46 necessitate coordination between multiple manipulators. During inference, this factor graph can
47 be treated as a probabilistic graphical model, where the learned skill factor and spatial constraint
48 factor distributions are composed to form a joint distribution of complete plans. Through 13 long-

49 horizon manipulation tasks in simulation and the real world, we show that GFC can solve complex
 50 bimanual manipulation tasks and exhibits strong generalization to unseen planning tasks with novel
 51 combinations of objects and constraints.

52 2 Background

53 **Diffusion Models.** A core component of our method is based on distributions learned using dif-
 54 fusion models. A diffusion model learns an unknown distribution $p(\mathbf{x}^{(0)})$ from its samples by ap-
 55 proximating the score function $\nabla \log p$. It consists of two processes: a *forward diffusion or noising*
 56 process that progressively injects noise and a *reverse diffusion or denoising* process that iteratively
 57 removes noise to recover clean data. The forward process simply adds Gaussian noise ϵ to clean
 58 data as $\mathbf{x}^{(t)} = \mathbf{x}^{(0)} + \sigma_t \epsilon$ for a monotonically increasing σ_t . The reverse process relies on the score
 59 function $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}^{(t)})$ where p_t is the distribution of noised data $\mathbf{x}^{(t)}$. In practice, the unknown
 60 score function is estimated using a neural network $\epsilon_\phi(\mathbf{x}^{(t)}, t)$ by minimizing the denoising score
 61 matching [9] objective $\mathbb{E}_{t, \epsilon, \mathbf{x}^{(0)}} [\lambda(t) \|\epsilon - \epsilon_\phi(\mathbf{x}^{(t)}, t)\|^2]$ where $\lambda(t)$ is a time-dependent weight. Sev-
 62 eral recent works have explored the advantages of diffusion models like scalability [10, 11, 12, 13]
 63 and the ability to learn multi-modal distributions [14, 15, 16, 17]. We are particularly interested in
 64 the compositional ability [18, 19, 20, 21, 6] of these models for the proposed method.

65 **Problem setup.** We assume access to a library of parameterized skills [22] $\pi \sim \Pi$ such as primitive
 66 actions like `Pick` and `Place`. Each skill π requires a pre-condition to be fulfilled and is parame-
 67 terized by a continuous parameter $a \in A_\pi$ governing the desired motion while executing the skill
 68 in a state s . For a given symbolically feasible task plan from a starting state s_0 to reach a specified
 69 goal condition s_{goal} , generated by a task planner or given by an oracle, the problem is to obtain
 70 the sequence of continuous parameters to make the plan geometrically feasible. For example, given
 71 a nail at a target location and a hammer on a table, the symbolic plan is to `Pick` the hammer and
 72 `Reach` the nail. A geometrically-feasible plan requires suitable `Pick` and `Reach` parameters such
 73 that the hammer’s head can strike the nail.

74 **Learning for skill chaining.** Existing works along this direction model the planning problem as
 75 a “chaining” problem: They first model the pre-conditions and effect state distributions for every
 76 skill $\pi \sim \Pi$ from the available data and a symbolic *plan skeleton* $\Phi_K = \{\pi_1, \pi_2, \dots, \pi_K\}$ consisting
 77 of K -skills is constructed. With this model, they search for the given skill sequence (plan) such
 78 that each skill satisfies the pre-conditions of the next skill in the plan. STAP [5] used learned pri-
 79 ors to perform data-driven optimization with the cross-entropy maximization method. In GSC [6],
 80 the policy and transition model is formulated as a diffusion model based distribution $p_\pi(s, a_\pi, s')$
 81 which allows for flexible chaining. While the forward chain ensures dynamics consistency in the
 82 plan, backward chain ensures that the goal is reachable from the intermediate states. For a forward
 83 rollout trajectory $\tau = \{s_0, a_{\pi_1}, s_1, a_{\pi_2}, s_{goal}\}$ associated with skeleton $\Phi_2 = \{\pi_1, \pi_2\}$, the resulting
 84 forward-backward combination based on GSC [6] can be represented as

$$p_\tau(\tau | s_0, s_{goal}) \propto \frac{p_{\pi_1}(s_0, a_{\pi_1}, s_1) p_{\pi_2}(s_1, a_{\pi_2}, s_{goal})}{\sqrt{p_{\pi_1}(s_1) p_{\pi_2}(s_1)}} \quad (1)$$

85 3 Method

86 We aim to solve unseen long-horizon planning problems by exploiting the inter-dependencies be-
 87 tween the objects important for the task at hand in the scene. Our method uses a spatial-temporal
 88 factor graph [8] to represent states and realize their temporal evolution by the application of skills.
 89 While previous works have considered *vectorized state* representations making it difficult to de-
 90 couple spatial-independence, we focus on *factorized state* representations such that the state of the
 91 environment is entirely modular, containing information about all the objects in the scenario and
 92 the task-specific constraints between them. We transform the graph into a probabilistic graphical
 93 model by representing temporal factors as skill-level transition distributions and spatial factors as
 94 constraint-satisfaction distributions. A composition of all the factors jointly represents sequential
 95 and coordinated manipulation plans directly at inference and can be solved by sampling optimal
 96 node variables using reverse diffusion sampling.

97 **3.1 Representing States, Skills, and Plans in Factor Graphs**

98 **States as factor graphs.** We define a factor graph $\{\mathcal{V}, \mathcal{F}\}$ of a state s consisting of the decision
 99 variable \mathcal{V} and factor \mathcal{F} nodes. Every robot and object is represented as a decision variable node
 100 $v \in \mathcal{V}$ containing their respective state. Factors $f \in \mathcal{F}$ between nodes in a given state are *spatial*
 101 *constraints*. For example, a Grasped spatial factor specifies admissible rigid transforms between a
 102 gripper and an object. Mathematically, we construct a probabilistic graphical model from to formu-
 103 late the distribution of a state, $p(s)$ as the composition of all the factor distributions:

$$p(s) \propto \prod_{f \in \mathcal{F}} p_f(S_f) \quad \text{where } s \equiv \bigcup_{f \in \mathcal{F}} S_f \quad (2)$$

104 where $p_f(S_f)$ represents the joint factor potential of nodes $v \in S_f \subseteq \mathcal{V}$, i.e. all nodes involved in a
 105 factor ¹ and s is the joint distribution of all such nodes. This indicates that the joint distribution of
 106 all the nodes must satisfy each of the factors, also explored by Diffusion-CCSP [21].

107 **Skills as temporal factors.** To represent transitions between states, we adapt parameterized
 108 skills [22] for a factor graph formulation. We define the preconditions of a skill as a set of nodes
 109 and factors, thus considering a skill feasible *iff* the precondition factors are satisfied. For example,
 110 for state s_0 illustrated in Figure 1, the nodes of a factor graph are $\{L_0, H_0, R_0, B_0\}$ and the factors
 111 existing in this scene are $\{\text{Grasped}(L_0, H_0)=\text{True}\}$. Now, since this factor is a precondition of the
 112 skill $\text{Move}(L_0, H_0)$ that moves the hammer in hand to align with the box, it must be satisfied for the
 113 skill to be feasible. The effect of executing a skill creates a new factor graph s' by changing the state
 114 of the nodes involved and, optionally, adding or removing their factors. This results in a *temporal*
 115 *factor* between the transitioned nodes of s and s' with the continuous action parameter of the skill
 116 a_π . The skill definitions can be extracted from standard PDDL symbolic skill operator with minor
 117 adaptations, following the duality of factor graphs and plan skeletons [1]. Eventually, we solve an
 118 optimization problem: satisfying the Aligned, Grasped, and the transition dynamics constraints by
 119 finding the correct Move parameters a_{π_1} . Each skill in a plan introduces additional nodes and factors
 120 to the factor graph, with added complexity for optimization.

121 Mathematically, we can use the distribution $p(s)$ as established in Equation 2 with all the spatial
 122 factors, and represent the temporal skill factor distribution of k^{th} -skill π_k as the joint distribution:
 123 $p_{\pi_k}(s, a, s') \equiv p_{\pi_k}(S_{\pi_k}, a, S'_{\pi_k})$, $S_{\pi_k} \subseteq \mathcal{V}_{pre}^{\pi_k}$ which is executable *iff* the skill’s pre-condition
 124 $S_{pre}^{\pi_k} \equiv \{\mathcal{V}_{pre}^{\pi_k}, \mathcal{F}_{pre}^{\pi_k}\}$ is satisfied by the current state i.e. $\mathcal{V}_{pre}^{\pi_k} \subseteq \mathcal{V}$ and $\mathcal{F}_{pre}^{\pi_k} \subseteq \mathcal{F}$. Once executed,
 125 it leads to the transitioned state S'_{π_k} . Based on the above formulation of a short-horizon transition
 126 distribution, we extend to construct a plan-level distribution as already established by GSC [6] and
 127 shown in Equation 1. We leverage the modularity of factored states by replacing states s with a set
 128 of decision variables S_{π_k} in the interest of skill π_k allowing us to rewrite Equation 1 as:

$$p(\tau) \propto \frac{\prod_{\pi_k \in \Phi} p_{\pi_k}(v_k \in \mathcal{V}_{pre}^{\pi_k}, a_k, v'_k \in \mathcal{V}_{effect}^{\pi_k})}{\sqrt{\prod_{v_i \in \mathcal{V}_i} p_{\pi_{i-}}(v_i) p_{\pi_{i+}}(v_i)}} \quad (3)$$

129 if we consider that some set of intermediate nodes \mathcal{V}_i are connected by two sequential skills π_{i-} and
 130 π_{i+} i.e. $\mathcal{V}_i \in S'_{\pi_{i-}} \cap S_{\pi_{i+}}$.

131 **Representing coordination.** A key advantage of the factor graph representation is the ability to
 132 model multi-arm coordination tasks by connecting the temporal chains of each arm using spatial
 133 constraints. Such tasks often require skills to be simultaneously executed on each arm to operate
 134 on different or the same objects. We consider two cases for parallel skill execution, where multiple
 135 robots are operating on: (1) independent objects and (2) the same object, leading to independent and
 136 dependent temporal chains respectively. With our factorized state representation, we can indepen-
 137 dently control the execution of individual skills correlated with the nodes of interest and calculate
 138 the cumulative effect by applying the union of the effects of all the skills to the current factor graph.
 139 We consider a scenario shown in Figure 2 (Left). The left and right gripper arm L_0, R_0 are hold-
 140 ing the pink C_0 and green M_0 cup ($\{\text{Grasped}(L_0, C_0)=\text{True}\}$ and $\{\text{Grasped}(R_0, M_0)=\text{True}\}$)

¹i.e. a factor f is included *iff* there is an edge between f and some $v \in \mathcal{V}$ which also implies $v \in S_f \subseteq \mathcal{V}$.

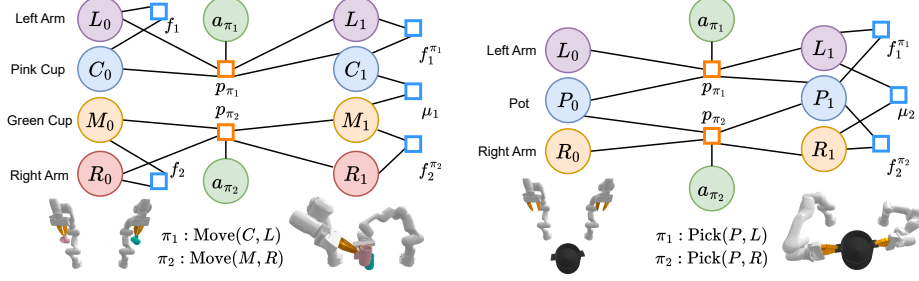


Figure 2: (Left) **Parallel independent chaining** The figure shows the execution of two skills (π_1 and π_2) in-parallel on two independent sets of nodes (L, C and R, M) to modify their existing factors (Grasped). The two independent executions can be connected via external factors μ_1 (FixedTransform) introducing spatial dependencies between nodes C and M. (Right) **Parallel dependent chaining** The figure shows overlapping nodes of interest while parallel execution of two skills. The pot is to be picked by using both arms simultaneously. The effect of this is resulting factors (Grasped) between (L, P and R, P) and external factor μ_2 (FixedTransform) between L and R. Overlapping nodes satisfy both skill’s temporal effects.

141 respectively. While both the grippers can independently execute the skill Move to modify separate
 142 factors ($f_1^{\pi_1}$ and $f_2^{\pi_2}$), one can add a constrained relationship factor (μ_1) between the two mugs
 143 representing a set of transforms that satisfy the precondition of Pour. Such an ability to augment
 144 constraints flexibly allows zero-shot coordination planning for unseen tasks at test time even with
 145 parallel skill executions on the same object as shown in Figure 2 (Right).

146 3.2 Generative Factor Chaining

147 Now we have a formulation to construct a symbolic spatial-temporal factor graph plan for a task and
 148 chain them using spatial factor and temporal skill factors sequentially or in parallel. To make this
 149 plan geometrically feasible, we must find the optimal node variable values. While classical solvers
 150 require modeling the transition dynamics of complex manipulation tasks, sampling-driven optimiza-
 151 tion with learned models provides less flexibility and modularity [6]. In this work, we leverage the
 152 expressive generative model to capture the transition dynamics and exploit the compositionality of
 153 diffusion models. Given a symbolically feasible factor graph plan, our method, termed Generative
 154 Factor Chaining (GFC), can flexibly compose spatial-temporal factor distributions to sample optimal
 155 node variable values for the complete plan.

156 **Probabilistic model for trajectory plan as spatial-temporal factor graphs.** Now, we again con-
 157 sider the spatial graph for representing the state, where the probability of finding a state s is the joint
 158 distribution of all the nodes in the factor graph. We will now integrate the spatial factors with the
 159 temporal factors considering the compensation term introduced in Equation 2 and Equation 3 along
 160 with the constraint factors across the chain $\mu \in \mathcal{M}$ as:

$$p(\tau) \propto \frac{\prod_{\pi_k \in \Phi} p_{\pi_k}(v_k \in \mathcal{V}_{pre}^{\pi_k}, a_k, v_{k+1} \in \mathcal{V}_{effect}^{\pi_k}) \prod_{k=0}^K \prod_{f \in \mathcal{F}_k} p_f(\mathcal{S}_f)}{\sqrt{\prod_{v_i \in \mathcal{V}_i} p_{\pi_{i-}}(v_i) p_{\pi_{i+}}(v_i)}} \prod_{\mathcal{M}} f_{\mu}(S_{\mu}) \quad (4)$$

161 This completes the joint distribution of all the nodes in the spatial-temporal factor graph plan consid-
 162 ering the temporal factors for all skills with their pre-condition and effect nodes, all spatial factors for
 163 all states in the plan, and all intermediate nodes in the temporal chain. We show our implementation
 164 of this formulation in algorithm 1.

165 For the sake of simplicity, we will formulate the probabilistic model for the two chains shown in Fig-
 166 ure 2 by following the forward-backward analysis introduced by GSC and discussed in section 2.
 167 We can write the bottom chain can be constructed based on Equation 4 as:

$$\frac{p_{\pi_1}(L_0, P_0, a_{\pi_1}, L_1, P_1) p_{\pi_2}(R_0, P_0, a_{\pi_2}, R_1, P_1)}{\sqrt{p_{\pi_1}(P_1) p_{\pi_2}(P_1)}} p_{\mu_2}(L_1, R_1) \quad (5)$$

168 where the factors are dependent on each other. It is worth noting that the augmented constraint
 169 factors p_{μ} work as a weighing function and can be more precisely represented by $p_{\mu}(S_{\mu}) \equiv p_{\mu}(y =$
 170 $1 | S_{\mu})$ for some constraint-satisfaction index y .

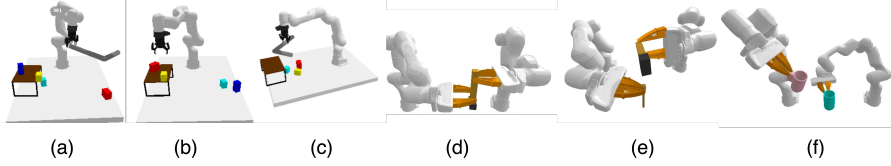


Figure 3: Evaluation tasks: **(a) Hook reach:** Hook is used to pull an object in the robot’s workspace followed by other skills. **(b) Constrained packing:** Multiple objects must be placed on a rack without collisions. **(c) Rearrangement push:** Hook is used to push objects to a desired arrangement followed by other skills. **(d) Hammer place:** A hammer must be handed over to another manipulator and placed in a target box. **(e) Hammer nail:** A hammer must be handed over to another manipulator and a configuration must be achieved to strike a nail. **(f) Pour cup:** Cups must be brought in a configuration that allows successful pouring from one to another.

171 We align towards diffusion model-based learned distributions to represent the probabilities in the
 172 formulated probabilistic graphical model. We transform the probabilities into their respective score
 173 functions $\epsilon(\mathbf{x}^{(t)}, t)$ for a particular reverse diffusion sampling step t and train it using score matching
 174 loss. Hence, for sampling a scene-graph for Equation 4, we have

$$\begin{aligned} \epsilon(\tau^{(t)}, t) = & \sum_{\pi_k \in \Phi} \epsilon_{\pi_k}(v_k^{(t)} \in \mathcal{V}_{pre}^{\pi_k}, a_k^{(t)}, v_{k+1}^{(t)} \in \mathcal{V}_{effect}^{\pi_k}, t) + \sum_{k=0}^K \sum_{f \in \mathcal{F}_k} \epsilon_f(\mathcal{S}_f^{(t)}, t) \\ & - \frac{1}{2} \sum_{v_i \in \mathcal{V}_i} \left[\epsilon_{\pi_{i-}}(v_i^{(t)}, t) \epsilon_{\pi_{i+}}(v_i^{(t)}, t) \right] + \sum_{\mathcal{M}} \epsilon_{f_\mu}(S_\mu^{(t)}, t) \end{aligned}$$

175 Following this, we can show for the dependent factor chain in Equation 5 as:

$$\begin{aligned} \epsilon(L_0^{(t)}, P_0^{(t)}, R_0^{(t)}, L_1^{(t)}, P_1^{(t)}, R_1^{(t)}, t) = & \epsilon_{\pi_1}(L_0^{(t)}, P_0^{(t)}, a_{\pi_1}^{(t)}, L_1^{(t)}, P_1^{(t)}, t) + \\ \epsilon_{\pi_2}(R_0^{(t)}, P_0^{(t)}, a_{\pi_2}^{(t)}, R_1^{(t)}, P_1^{(t)}, t) - & \frac{1}{2} \epsilon_{\pi_1}(P_1^{(t)}, t) - \frac{1}{2} \epsilon_{\pi_2}(P_1^{(t)}, t) + \epsilon_{\mu_2}(L_1^{(t)}, R_1^{(t)}, t) \end{aligned}$$

176 Such a representation leads to a cumulative score calculation of the joint distribution of all the nodes
 177 of interest to the factor using linear addition and subtraction. We can realize from Equation 3.2
 178 that the final score function depends on the composition of all the factors in the spatial-temporal
 179 factor graph. While factors $f \in \mathcal{F}$ are mostly modeled implicitly by the temporal skills, the external
 180 factors can be any arbitrary spatial constraints that ensure the satisfaction of the pre-condition of the
 181 subsequent skills. Hence, with new additions to the set of external factors $\mu' \in \mathcal{M}'$, one can reuse
 182 the same temporal skills with added new spatial constraints.

183 4 Experiment

184 In this section, we seek to validate the following hypotheses: (1) GFC relaxes strict temporal depen-
 185 dency to allow spatial-temporal reasoning, performing better or on par with prior works in single-
 186 arm long-horizon sequential manipulation tasks, (2) GFC can effectively solve unseen coordination
 187 tasks, and (3) GFC is adept in reasoning about long-horizon action dependency while being robust
 188 to increasing task horizons. We systematically evaluated our method on 9 long-horizon single-arm
 189 manipulation tasks from prior works and 4 complex multi-arm coordination tasks in simulation. We
 190 also demonstrate deploying GFC on a bimanual Franka Panda setup in the real world.

191 **Relevant baselines and metrics:** Our proposed method is based on factorized states and supports
 192 long-horizon planning for collaborative tasks directly at inference via probabilistic chaining. In this
 193 context, we consider prior methods based on probabilistic chaining with vectorized states (**GSC** [6])
 194 and discriminative search-based approaches for solving long-horizon planning by skill chaining:
 195 with uniform priors (**Random CEM** or **RCEM**) or learned policy priors (**STAP** [5]). Since all
 196 prior works use sequential planning, we compare the performance of the proposed method on the
 197 sequential version of the parallel skeleton. Further information on data collection, training of skill
 198 diffusion models and real robot setup is provided in [Supp. S3](#) and [Supp. S4](#) respectively.

199 **GFC relaxes strict linear dependency assumptions.** We first evaluate GFC on single-manipulator
 200 long-horizon tasks introduced by STAP [5]. These tasks consider manipulation by reasoning about
 201 the usage of a tool (a hook) to manipulate blocks out of or into the robot workspace (sample initial
 202 states shown in Figure 3(a-c)) and provide the descriptions of each considered task in the caption.

Table 1: We show performance comparison of our method with relevant baselines on 9 single manipulator tasks and 3 two-manipulator tasks based on 100 trials for each of them. The task length shows the relative difficulty of solving them. We also conduct evaluation on 3 extended tasks to show robustness of GFC to task length ($|\mathcal{T}|$) and efficient reasoning about interstep dependencies.

Evaluation Tasks		RCEM	DAF [4]	STAP [5]	GSC [6]	GFC	$ \mathcal{T} $	
Single Manipulator	Hook Reach	T1	0.54	0.32	0.88	0.84	0.82	4
		T2	0.40	0.05	0.82	0.84	0.82	5
		T3	0.30	0.00	0.76	0.76	0.80	5
	Rearrangement Push	T1	0.30	0.0	0.40	0.68	0.68	4
		T2	0.10	0.08	0.52	0.60	0.65	6
		T3	0.02	0.0	0.18	0.18	0.25	8
	Constrained Packing	T1	0.45	0.45	0.65	0.75	0.75	6
		T2	0.45	0.70	0.68	1.0	1.0	6
		T3	0.10	0.0	0.20	1.0	1.0	8
Bimanual Manipulation	Hammer Place	0.05	-	0.28	0.41	0.63	8	
	Pour Cup	0.10	-	0.18	0.15	0.41	4	
	Hammer Nail	0.02	-	0.15	0.15	0.34	11	
Longer Horizon Evaluation Tasks								
Handback Hammer Nail						0.24	16	
Handback Hammer Nail w/ auxilliary tasks						0.25	18	
Handback Hammer Nail w/ extended auxilliary tasks						0.21	20	

203 While these tasks are originally designed to highlight linear sequential dependencies, there are steps
 204 with indirect dependencies or independence that only GFC can effectively model because of the
 205 factorized states. For example, in *Rearrangement Push*, the picking pose of the cube should not
 206 affect the tool use steps. As shown in Table 1, we observe that the performance of GFC is con-
 207 sistent on-par with the baseline for tasks with strict linear dependencies such as *Hook Reach* and
 208 on-par or better for tasks with more complex dependency structures such as *Rearrangement Push*.
 209 This validates our hypothesis that GFC effectively models sequential dependencies, in addition to
 210 independence and skipped-step dependencies in long-horizon tasks.

211 **GFC can solve complex coordinated manipulation tasks.** Here,
 212 we aim to validate that GFC can effectively plan and solve different
 213 types of coordinated manipulation tasks. We present results on tasks
 214 with increased collaboration challenges. First, we consider tasks
 215 that require coordination but can be serialized into interleaved skill
 216 chains and solved by prior skill-chaining methods. *Hammer Place*,
 217 as shown in Figure S16, is for one arm to pick a hammer, hand it
 218 over to another arm for placement into a target box. *Hammer*
 219 *Nail* is an extension where, after hammer handover, first arm picks
 220 up a nail and both arms coordinate to move to positions such that the
 221 hammer’s head is aligned with the nail for the subsequent striking
 222 step. The task is illustrated in Figure S16. As evident from Table 1,
 223 GFC outperforms all baselines in both tasks. The gap is larger in
 224 the more challenging *Hammer Nail* task, which includes additional
 225 spatial and temporal constraints such as the hammer must be re-
 226 grasped towards the tail end for the subsequent hammering step,
 227 and the hammer and nail must be aligned for a successful strike.
 228 This demonstrates that GFC can effectively model and resolve both spatial and temporal constraints
 229 in complex tasks.

230 **GFC can zero-shot generalize to new bimanual tasks by composing single-arm skill chains.**
 231 The *Pour Cup* (Figure S11) task is to Pick a cup with each arm, Move to position the two cups, and
 232 Pour the content of one into the other. GFC can directly reuse Pick and Move skill models and adapt
 233 the Strike skill model for the Pour step by adding a new spatial constraint. Unlike hammer that can
 234 strike from either face of the head, the cups can only be poured using the open top and not the closed
 235 bottom. The constraint can be directly added as a spatial factor. A quantitative comparison is shown
 236 in Table 1. Finally, we consider the *Bimanual Reorientation* (Figure S12) task where two arms
 237 must simultaneously operate on the same object of interest (a pot), lift it up, and rotate it to a target
 238 reorientation angle (about z-axis) as illustrated in Figure 4 (Top) for a 45-deg angle. The tasks must
 239 be solved via parallel skill chaining with spatial constraints and hence none of the prior baselines can
 240 be used. The factor graph (Figure 2 Right) includes a spatial fixed transform constraint between both
 241 the arms and hence the subsequent skills operate while satisfying the constraint. Figure 4 (Bottom)

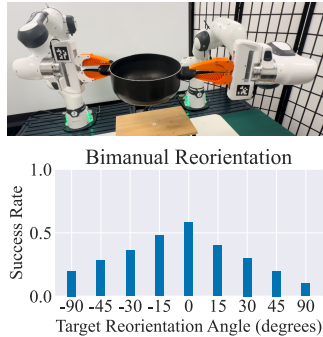


Figure 4: Evaluating GFC on bimanual reorientation where two arms simultaneously pick and reorient a pot.

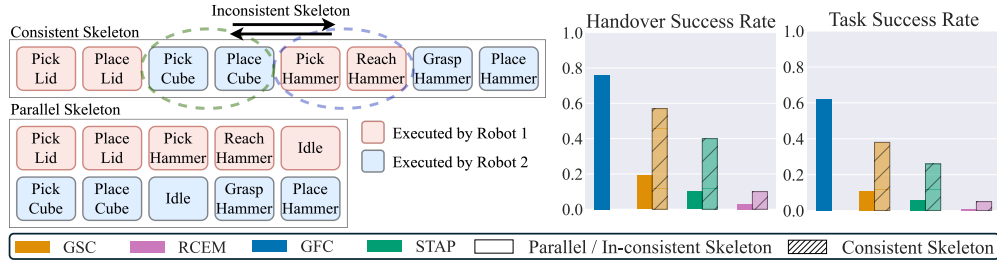


Figure 5: **Linear chaining has limitations.** Baseline methods with linear chain assumption suffers from performance drop when given inconsistent skill chains, where steps with sequential dependencies are swapped. GFC retains high success rate using the parallel skeleton.

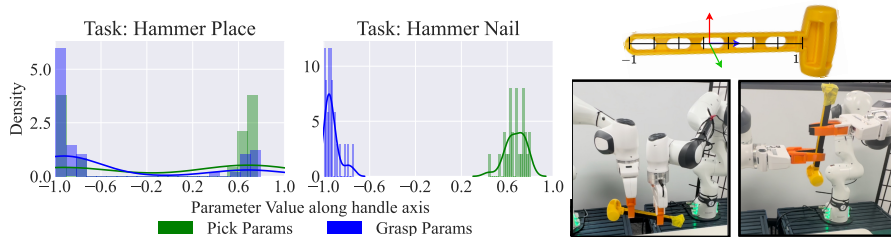


Figure 6: **Analysis of coordination.** We show that the planner is able to reason about the long-horizon action dependency of Pick and Grasp skills. (Left) While we see that *Hammer Place* can be solved by pick/grasp at head/tail and vice versa, to satisfy the precondition of Strike in *Hammer Nail*, the hammer must be grasped near tail so must be picked near head. (Right) We show orientation reasoning, where the hammer can either be grasped on the same side or the flip side.

242 shows a detailed task success rate breakdown given different orientation goals. The spatial and
 243 temporal challenges posed by the task are detailed further in [Supp. S5](#).

244 **GFC can handle independence and inconsistent skill chains.** Here, we analyze how *independent*
 245 steps in a sequential manipulation chain affects the performance of each method. We consider *Ham-*
 246 *mer Place*, where the order of transporting the cube and handing over hammer is interchangeable.
 247 As illustrated in [Figure 5](#), we consider a *consistent* plan skeleton where sequentially-dependent steps
 248 for the two main objectives, i.e., (1) opening lid then transporting cube and (2) picking, handing over,
 249 and placing hammers, are completely sequentially. We also consider an *inconsistent* plan skeleton
 250 where the steps are interleaved. We show the handover success and overall task success in [Fig-](#)
 251 [ure 5](#) (Right). A successful handover requires choosing compatible parameters for Pick, Regrasp,
 252 and Move skills. While this increases the difficulty leading to lower scores in the handover suc-
 253 cess rate, even with a minor distraction in *inconsistent* skeleton, the previous approaches failed to
 254 propagate the skipped-step dependencies as evident from the task success rate.

255 **GFC can reason about action dependency while being robust to increasing task horizons.** We
 256 observe in [Figure 6](#) (left) that while *Hammer Place* task can be solved by picking or grasping on
 257 any end of the hammer handle, *Hammer Nail* requires more constrained parameter sampling. Fur-
 258 ther, in addition to the parameter selection along the handle axis, the method also samples suitable
 259 orientation (same or flip side) for grasping as shown by two examples in [Figure 6](#) (right). We further
 260 give an example of the capability of our method in handling longer horizon inter-step dependencies
 261 in [Figure S17](#) and simultaneously being robust with respect to the task length as shown in [Table 1](#).

262 5 Conclusion

263 We presented GFC, a learning-to-plan method for complex coordinated manipulation tasks. GFC
 264 can flexibly represent multi-arm manipulation with one or more objects with a spatial-temporal
 265 factor graph. During inference, GFC composes factor graphs where each factor is a diffusion model
 266 and samples long-horizon plans with reverse denoising. GFC is shown to solve sequential and
 267 coordinated tasks directly at inference and reason about long-horizon action dependency across
 268 multiple temporal chains. Our framework generalizes well to unseen multiple-manipulator tasks.

References

- 269
- 270 [1] C. R. Garrett, R. Chitnis, R. Holladay, B. Kim, T. Silver, L. P. Kaelbling, and T. Lozano-Pérez.
271 Integrated task and motion planning. *Annual review of control, robotics, and autonomous*
272 *systems*, 4:265–293, 2021.
- 273 [2] P.-L. Bacon, J. Harb, and D. Precup. The option-critic architecture. In *Proceedings of the AAAI*
274 *conference on artificial intelligence*, volume 31, 2017.
- 275 [3] D. Driess, J.-S. Ha, and M. Toussaint. Learning to solve sequential physical reasoning prob-
276 lems from a scene image. *The International Journal of Robotics Research*, 40(12-14):1435–
277 1466, 2021.
- 278 [4] D. Xu, A. Mandlekar, R. Martín-Martín, Y. Zhu, S. Savarese, and L. Fei-Fei. Deep affor-
279 dance foresight: Planning through what can be done in the future. In *2021 IEEE International*
280 *Conference on Robotics and Automation (ICRA)*, pages 6206–6213. IEEE, 2021.
- 281 [5] C. Agia, T. Migimatsu, J. Wu, and J. Bohg. Taps: Task-agnostic policy sequencing. *arXiv*
282 *preprint arXiv:2210.12250*, 2022.
- 283 [6] U. A. Mishra, S. Xue, Y. Chen, and D. Xu. Generative skill chaining: Long-horizon skill
284 planning with diffusion models. In *7th Annual Conference on Robot Learning*, 2023. URL
285 <https://openreview.net/forum?id=HtJE9ly5dT>.
- 286 [7] J. Liang, M. Sharma, A. LaGrassa, S. Vats, S. Saxena, and O. Kroemer. Search-based task
287 planning with learned skill effect models for lifelong robotic manipulation. In *2022 Interna-*
288 *tional Conference on Robotics and Automation (ICRA)*, pages 6351–6357. IEEE, 2022.
- 289 [8] F. Dellaert. Factor graphs: Exploiting structure in robotics. *Annual Review of Control,*
290 *Robotics, and Autonomous Systems*, 4:141–166, 2021.
- 291 [9] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based gener-
292 ative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*,
293 2020.
- 294 [10] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural*
295 *Information Processing Systems*, 33:6840–6851, 2020.
- 296 [11] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *arXiv preprint*
297 *arXiv:2010.02502*, 2020.
- 298 [12] Q. Zhang, M. Tao, and Y. Chen. gddim: Generalized denoising diffusion implicit models.
299 *arXiv preprint arXiv:2206.05564*, 2022.
- 300 [13] Q. Zhang and Y. Chen. Fast sampling of diffusion models with exponential integrator. *arXiv*
301 *preprint arXiv:2204.13902*, 2022.
- 302 [14] J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
303 2022.
- 304 [15] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *Advances in*
305 *Neural Information Processing Systems*, 34:8780–8794, 2021.
- 306 [16] A. Lugmayr, M. Danelljan, A. Romero, F. Yu, R. Timofte, and L. Van Gool. Repaint: In-
307 painting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF*
308 *Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022.
- 309 [17] A. Ajay, Y. Du, A. Gupta, J. Tenenbaum, T. Jaakkola, and P. Agrawal. Is conditional generative
310 modeling all you need for decision-making? *arXiv preprint arXiv:2211.15657*, 2022.
- 311 [18] Q. Zhang, J. Song, X. Huang, Y. Chen, and M.-Y. Liu. Diffcollage: Parallel generation of large
312 content with diffusion models. *ArXiv*, abs/2303.17076, 2023.
- 313 [19] T. Yu, T. Xiao, A. Stone, J. Tompson, A. Brohan, S. Wang, J. Singh, C. Tan, J. Peralta,
314 B. Ichter, et al. Scaling robot learning with semantically imagined experience. *arXiv preprint*
315 *arXiv:2302.11550*, 2023.

- 316 [20] Y. Du, S. Li, and I. Mordatch. Compositional visual generation with energy based models.
317 *Advances in Neural Information Processing Systems*, 33:6637–6647, 2020.
- 318 [21] Z. Yang, J. Mao, Y. Du, J. Wu, J. B. Tenenbaum, T. Lozano-Pérez, and L. P. Kaelbling. Com-
319 positional diffusion-based continuous constraint solvers. *arXiv preprint arXiv:2309.00966*,
320 2023.
- 321 [22] L. P. Kaelbling and T. Lozano-Pérez. Learning composable models of parameterized skills.
322 In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 886–893.
323 IEEE, 2017.
- 324 [23] R. S. Sutton, D. Precup, and S. Singh. Between mdps and semi-mdps: A framework for tem-
325 poral abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- 326 [24] D. Shah, P. Xu, Y. Lu, T. Xiao, A. Toshev, S. Levine, and B. Ichter. Value function spaces:
327 Skill-centric state abstractions for long-horizon reasoning. *arXiv preprint arXiv:2111.03189*,
328 2021.
- 329 [25] O. Nachum, S. S. Gu, H. Lee, and S. Levine. Data-efficient hierarchical reinforcement learning.
330 *Advances in neural information processing systems*, 31, 2018.
- 331 [26] W. Masson, P. Ranchod, and G. Konidaris. Reinforcement learning with parameterized actions.
332 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016.
- 333 [27] I. Kapelyukh, V. Vosylius, and E. Johns. Dall-e-bot: Introducing web-scale diffusion models
334 to robotics. *IEEE Robotics and Automation Letters*, 2023.
- 335 [28] T. Pearce, T. Rashid, A. Kanervisto, D. Bignell, M. Sun, R. Georgescu, S. V. Macua, S. Z. Tan,
336 I. Momennejad, K. Hofmann, et al. Imitating human behaviour with diffusion models. *arXiv*
337 *preprint arXiv:2301.10677*, 2023.
- 338 [29] C. Chi, S. Feng, Y. Du, Z. Xu, E. A. Cousineau, B. Burchfiel, and S. Song. Diffusion policy:
339 Visuomotor policy learning via action diffusion. *ArXiv*, abs/2303.04137, 2023.
- 340 [30] U. A. Mishra and Y. Chen. Reorientdiff: Diffusion model based reorientation for object ma-
341 nipulation. *arXiv preprint arXiv:2303.12700*, 2023.
- 342 [31] Y. Du, M. Yang, B. Dai, H. Dai, O. Nachum, J. B. Tenenbaum, D. Schuurmans, and P. Abbeel.
343 Learning universal policies via text-guided video generation. *ArXiv*, abs/2302.00111, 2023.
- 344 [32] M. Reuss, M. Li, X. Jia, and R. Lioutikov. Goal-conditioned imitation learning using score-
345 based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.
- 346 [33] M. Reuss and R. Lioutikov. Multimodal diffusion transformer for learning from play. In *2nd*
347 *Workshop on Language and Robot Learning: Language as Grounding*, 2023. URL <https://openreview.net/forum?id=nvtxqMGpn1>.
- 349 [34] M. Janner, Y. Du, J. Tenenbaum, and S. Levine. Planning with diffusion for flexible behavior
350 synthesis. In *International Conference on Machine Learning*, 2022.
- 351 [35] J. Chen, J. Li, Y. Huang, C. Garrett, D. Sun, C. Fan, A. Hofmann, C. Mueller, S. Koenig, and
352 B. C. Williams. Cooperative task and motion planning for multi-arm assembly systems. *arXiv*
353 *preprint arXiv:2203.02475*, 2022.
- 354 [36] L. Nägele, A. Hoffmann, A. Schierl, and W. Reif. Legobot: Automated planning for coordi-
355 nated multi-robot assembly of lego structures. In *2020 IEEE/RSJ International Conference on*
356 *Intelligent Robots and Systems (IROS)*, pages 9088–9095. IEEE, 2020.
- 357 [37] L. P. Ureche and A. Billard. Constraints extraction from asymmetrical bimanual tasks and their
358 use in coordinated behavior. *Robotics and autonomous systems*, 103:222–235, 2018.
- 359 [38] F. Amadio, A. Colomé, and C. Torras. Exploiting symmetries in reinforcement learning of
360 bimanual robotic tasks. *IEEE Robotics and Automation Letters*, 4(2):1838–1845, 2019.

- 361 [39] R. Chitnis, S. Tulsiani, S. Gupta, and A. Gupta. Efficient bimanual manipulation using learned
362 task schemas. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*,
363 pages 1149–1155. IEEE, 2020.
- 364 [40] H. Ha, J. Xu, and S. Song. Learning a decentralized multi-arm motion planner. *arXiv preprint*
365 *arXiv:2011.02608*, 2020.
- 366 [41] J. Grannen, Y. Wu, B. Vu, and D. Sadigh. Stabilize to act: Learning to coordinate for bimanual
367 manipulation. In *Conference on Robot Learning*, pages 563–576. PMLR, 2023.
- 368 [42] A. Tung, J. Wong, A. Mandlekar, R. Martín-Martín, Y. Zhu, L. Fei-Fei, and S. Savarese. Learn-
369 ing multi-arm manipulation through collaborative teleoperation. In *2021 IEEE International*
370 *Conference on Robotics and Automation (ICRA)*, pages 9212–9219. IEEE, 2021.
- 371 [43] F. Dellaert. Factor graphs: Exploiting structure in robotics. *Annu. Rev. Control. Robotics*
372 *Auton. Syst.*, 4:141–166, 2021. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:234254791)
373 [234254791](https://api.semanticscholar.org/CorpusID:234254791).
- 374 [44] M. Toussaint. Logic-geometric programming: An optimization-based approach to combined
375 task and motion planning. In *IJCAI*, pages 1930–1936, 2015.
- 376 [45] Y. Lee, P. Huang, K. M. Jatavallabhula, A. Z. Li, F. Damken, E. Heiden, K. Smith,
377 D. Nowrouzezahrai, F. Ramos, and F. Shkurti. Stamp: Differentiable task and motion planning
378 via stein variational gradient descent. *arXiv preprint arXiv:2310.01775*, 2023.
- 379 [46] F. Dellaert. Factor graphs: Exploiting structure in robotics. *Annual Re-*
380 *view of Control, Robotics, and Autonomous Systems*, 4(1):141–166, 2021. doi:
381 [10.1146/annurev-control-061520-010504](https://doi.org/10.1146/annurev-control-061520-010504). URL [https://doi.org/10.1146/](https://doi.org/10.1146/annurev-control-061520-010504)
382 [annurev-control-061520-010504](https://doi.org/10.1146/annurev-control-061520-010504).
- 383 [47] Q. Xiao, Z. Zaidi, and M. Gombolay. Multi-camera asynchronous ball localization and trajec-
384 tory prediction with factor graphs and human poses. *arXiv preprint arXiv:2401.17185*, 2024.
- 385 [48] Y. Hao, Y. Gan, B. Yu, Q. Liu, S.-S. Liu, and Y. Zhu. Blitzcrank: Factor graph accelerator for
386 motion planning. In *2023 60th ACM/IEEE Design Automation Conference (DAC)*, pages 1–6.
387 IEEE, 2023.
- 388 [49] Z. Wang, C. R. Garrett, L. P. Kaelbling, and T. Lozano-Pérez. Active model learning and
389 diverse action sampling for task and motion planning. In *2018 IEEE/RSJ International Con-*
390 *ference on Intelligent Robots and Systems (IROS)*, pages 4107–4114. IEEE, 2018.
- 391 [50] Z. Wang, C. R. Garrett, L. P. Kaelbling, and T. Lozano-Pérez. Learning compositional models
392 of robot skills for task and motion planning. *The International Journal of Robotics Research*,
393 40(6-7):866–894, 2021.
- 394 [51] B. Kim, L. P. Kaelbling, and T. Lozano-Perez. Guiding the search in continuous state-action
395 spaces by learning an action sampling distribution from off-target samples. *arXiv preprint*
396 *arXiv:1711.01391*, 2017.
- 397 [52] X. Fang, C. R. Garrett, C. Eppner, T. Lozano-Pérez, L. P. Kaelbling, and D. Fox. Dimsam:
398 Diffusion models as samplers for task and motion planning under partial observability. *arXiv*
399 *preprint arXiv:2306.13196*, 2023.
- 400 [53] W. Peebles and S. Xie. Scalable diffusion models with transformers. *arXiv preprint*
401 *arXiv:2212.09748*, 2022.
- 402 [54] W. Crooks, G. Vukasin, M. O’Sullivan, W. Messner, and C. Rogers. Fin ray® effect inspired
403 soft robotic gripper: From the roboSoft grand challenge toward optimization. *Frontiers in*
404 *Robotics and AI*, 3:70, 2016.
- 405 [55] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead,
406 A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick. Segment anything. *arXiv:2304.02643*, 2023.

- 407 [56] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
408 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervi-
409 sion. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 410 [57] frankx. <https://github.com/pantor/frankx>, 2021.

411 S1 Main Contributions

412 **Generative Factor Chaining (GFC)** is proposed with the motivation of zero-shot motion planning
413 for long-horizon tasks. The goal is to use short-horizon skill transition distributions and efficiently
414 compose them to structure a long-horizon task-level distribution at inference. The factorized state
415 representation of GFC allows explicit reasoning of inter-object and skill-object interactions and
416 satisfying spatial constraints for coordinate manipulation. The primary contributions of GFC are as
417 follows:

- 418 1. A **generalized task representation** to formulate complex long-horizon coordination tasks
419 as a spatial-temporal factor graph of single-arm manipulation skill sequences connected via
420 spatial dependencies.
- 421 2. A **compositional framework** to compose short-horizon skill-level transition distributions
422 learned via diffusion models to represent long-horizon task-level distributions.
- 423 3. **Easy plug-and-play** via learning skill distributions with skill-level data only and add it
424 to the skill library. Any skill from the library can be plugged as temporal factors in the
425 spatial-temporal factor graph directly at inference for a given long-horizon task.

426 S2 Related Works

427 **Task and Motion Planning (TAMP).** TAMP frameworks decompose a complex planning problem
428 into constraint satisfaction problems at task and motion levels [23, 2, 24, 25, 26]. Notably, Garret
429 et al. [1] drew connections between TAMP and factor graphs [8], representing constraints as factors
430 and objects/robots as nodes. This formalism naturally allows reusing per-constraint solvers across
431 tasks. However, classical TAMP approaches often rely on accurate perception and system dynamics,
432 limiting their practical applications and scalability. We instead opt for a learning approach, while our
433 compositional factor graph representation remains heavily inspired by the classical TAMP paradigm.

434 **Generative models for planning.** Modern generative models have been applied to offline imitation
435 [19, 27, 28, 29, 30, 31, 32, 33] and reinforcement learning [34, 17]. In addition to modeling
436 complex state and action distributions, generative models have also been shown to encourage com-
437 positional generalization [18, 6, 20] by combining data across tasks [17, 34]. Most relevant to us are
438 Generative Skill Chaining (GSC) [6] and Diffusion-CCSP [21], both designed to achieve system-
439 atic compositional generalization. GSC composes skill chains through a guided diffusion process.
440 However, similar to other skill-chaining methods [4, 5], GSC cannot model non-linear dependencies
441 such as parallel skills and independence among skills. Diffusion-CCSP trains diffusion models to
442 generate object configurations to satisfy spatial constraints, while relying on external solvers to plan
443 the manipulation sequence. Our method is a unified framework to solve the combined problem: it
444 generates skill plans to satisfy both spatial and temporal constraints represented in a factor graph.

445 **Learning for coordinated manipulation.** Coordinating two or more arms for manipulation
446 presents numerous planning challenges [35, 36, 37], including the combinatorial search space
447 complex constraints for coordinated motion. Recent works have utilized learning-based frame-
448 works [38, 39, 40, 41, 42] in both Reinforcement Learning [38, 40] and offline Imitation Learn-
449 ing [42, 41]. However, most existing works have focused on learning task-specific policies [38, 41]
450 or require multi-arm demonstration data collected through a specialized teleoperation device [42].
451 In contrast, our factor graph-based representation enables solving multi-arm tasks by composing
452 multiple single-arm skills through inference-time optimization.

453 **Factor-graph representation for TAMP.** The graphical abstraction of a system for understand-
454 ing several inter-dependencies has been used in various domains [43]. Specifically in context to
455 task and motion planning (TAMP), such a representation allows the decomposition of multiple
456 modalities (discrete and continuous variables) in the state of a system [1]. Solving together for
457 discrete (logical decision variables) and continuous (motion parameters) can be formulated as a Hy-
458 brid Constraint Satisfaction Problem (H-CSP) problem, Logic-Geometric Program (LGP) [44], and
459 more recently by advanced gradient descent methods [45]. By following the factor-graph represen-
460 tation, the state space can be represented as a Cartesian product of all the subspaces and the action
461 space can be compactly represented based on the modalities they affect. We particularly follow the
462 *dynamic factor graph* representation used by Garrett et al. [1] to represent all the objects and action

463 parameters as the variable nodes of the graph and all the kinematic inter-dependencies as the factors
464 of the graph.

465 **Optimization for factor graphs.** Factor graphs are graphical models where the directed and undi-
466 rected factors, respectively represent the joint or conditional distribution of the variable nodes con-
467 nected to them. Most directed factors graphs as used for localization [46, 47] are formulated into
468 probabilistic graphical models of hidden-markov chains and solved for the maximum a posteri-
469 ori (MAP) [46, 48] estimates of the unknown node variables. Particularly in motion planning, opti-
470 mizing for all the variable nodes is often formulated as a constraint satisfaction problem [1, 21].

471 **Additional related works on learning for TAMP.** Recent works have shown that a number of
472 components of a TAMP system benefit from powerful generative models. Wang et al [49, 50] use
473 Gaussian Processes to learn continuous-space sampler for TAMP. Similarly, Kim et al. [51] use
474 GANs to learn action samplers. Fang et al. [52] propose to use Diffusion Models to capture complex
475 distributions such as Inverse Kinematics solutions, grasps, and contact dynamics. However, they
476 still rely on an overarching TAMP system to consume the generated samples to perform planning.
477 In contrast, our method directly forms a geometric plan sampler by chaining together factor-level
478 diffusion models.

Algorithm 1: Generative Factor Chaining (GFC) Algorithm

- 1 **Hyperparameters:**
 - 2 Number of reverse diffusion steps T
 - 3 **Inputs:**
 - 4 Pre-defined skill library $\Pi = \{\pi_1, \pi_2, \dots, \pi_M\}$
 - 5 Individual skill diffusion score functions ϵ_π
 - 6 Task skeleton $\Phi_K = \{\pi_0, \pi_1, \dots, \pi_K\}$: a sequence of skills of length K
 - 7 Scene graph sequence $\Phi_S = \{s_0, s_1, \dots, s_K\}$: a sequence of scene factors of length $K + 1$ where $s_k \equiv \{\mathcal{V}_k, \mathcal{F}_k\}$
 - 8 Goal condition $g \equiv \{\mathcal{V}_g, \mathcal{F}_g\}$
 - 9 Noise schedule σ
 - 10 Initialize $t = T = 1$
 - 11 Initialize Δt
 - 12 Initial node sequence $\mathbf{x}^{(T)} = \left[v_k^{(T)} \forall v \in \mathcal{V}_k, a_{\pi_k}^{(T)}, \dots \forall k \in [0, K] \right]$ sampled from $\mathcal{N}(\mathbf{0}, \sigma_T \mathbf{I})$
 - 13 **while** $t \geq 0$ **do**
 - 14 // Score of the joint distribution of all the nodes
 - 15 $\epsilon_\Phi(v_k^{(t)} \forall v \in \mathcal{V}_k, a_{\pi_k}^{(t)}, \dots \forall k \in [0, K], t) = \mathbf{0}$
 - 16 // Calculating the effective score of each node
 - 17 $\epsilon_\Phi(x^{(t)}, t) = \sum_{k=0}^K \epsilon_{\pi_k}(x^{(t)}, t) + \sum_{k=0}^K \sum_{f \in \mathcal{F}_k} \epsilon_f(x^{(t)}, t) \quad \forall x \in \mathbf{x}$ (Computational assumption, Equation 4)
 - 18 // Only for nodes connected with two temporal factors $f_{x,1}$ and $f_{x,2}$
 - 19 $\epsilon_\Phi(x^{(t)}, t) = \epsilon_\Phi(x^{(t)}, t) - \frac{1}{2} \left[\epsilon_{f_{x,1}}(x^{(t)}, t) + \epsilon_{f_{x,2}}(x^{(t)}, t) \right]$ (Denominator compensation, Equation 4)
 - 20 // calculating updated noised samples for the next reverse diffusion timestep
 - 21 $\tilde{\mathbf{x}}^{(t-1)} = \mathbf{x}^{(t)} + \dot{\sigma}_t \sigma_t \epsilon_\Phi(v_k^{(t)} \forall v \in \mathcal{V}_k, a_{\pi_k}^{(t)}, \dots \forall k \in [0, K], t) \Delta t$
 - 22 $t = t - \Delta t$
 - 23 **end**
 - 24 **Return** $\mathbf{x}^{(0)}$
-

479 S3 Experiment Setup, Model Training and Architecture

480 **Skill Data Collection and Skill Training** We consider a finite set of parameterized skills in our
481 skill library. While our framework supports flexible addition of new skills to the skill library, we
482 choose skills appropriate for the considered tasks. The parameterization, data collection, and train-
483 ing method for each of the skills is described as follows:

- 484 1. **Pick**: Gripper picks up an object from the table and the parameters contain 6-DoF pose in
485 the object’s frame of reference. The skill diffusion models are trained on successful pick
486 actions on all the available set of objects namely lid, cube, hammer, and nail/stake.
- 487 2. **Place**: Gripper places an object at the target location and parameters contain 6-DoF pose in
488 the place target’s frame of reference. This skill requires specifying two set of parameters,
489 the target pose and the target object (e.g. box, table). The picked object is placed and
490 successful placements are used to train the skill diffusion model.
- 491 3. **Move**: Gripper reaches a target location with an object in hand and parameters contain 6-
492 DoF pose in the manipulator’s frame of reference within the workspace. This skill captures
493 the distribution of the reachable workspace of the robot. When composed with the Move
494 skill of the second manipulator, the combined distribution captures the common workspace.
- 495 4. **ReGrasp**: Gripper grasps object mid-air and the parameters contain 6-DoF pose in the
496 object’s frame of reference. While collecting data directly for this skill is non-trivial, we
497 consider that if an object is picked up with parameters q_1 and moved with parameters q_2 ,
498 then the object can be grasped at the workspace location defined by q_2 with the ReGrasp
499 parameters as q_1 . Thus, we reuse Pick and Move data to train the skill diffusion model for
500 ReGrasp. While this is a design choice, with appropriate skill level data, we can train this
501 skill separately too.
- 502 5. **Push**: Gripper uses the grasped object to push away another object. The skill is motivated
503 from prior work [6, 5] where a hook object is used to Push blocks. The parameters of
504 this skill are (x, y, r, θ) such that the hook is placed at the (x, y) position on the table and
505 pushed by a distance r in the radial direction θ w.r.t. the origin of the manipulator. The
506 skill diffusion models is trained following GSC [6].
- 507 6. **Pull**: Gripper uses the grasped object to pull another object inwards. The skill is also mo-
508 tivated from prior work [6, 5] where a hook object is used to Pull blocks. The parameters
509 of this skill are (x, y, r, θ) such that the hook is placed at the (x, y) position on the table
510 and pulled by a distance r in the radial direction θ w.r.t. the origin of the manipulator. The
511 skill diffusion models is trained following GSC [6].
- 512 7. **Strike**: Gripper strikes another object with one object in hand (e.g., a hammer). As a
513 design choice, we do not train a skill diffusion model for this skill. Strike is primarily
514 used as a terminal skill. We are only concerned about the pre-condition as their effects can
515 be designed manually, which is similar to “subgoal skill” used in prior work. For example,
516 in order to satisfy the pre-condition of Strike, the hammer and nail must be aligned.
517 This can be satisfied in diverse configurations. However, the effect is achieved through a
518 deterministic motion.
- 519 8. **Pour**: Gripper rotates the object in hand in a pouring fashion. Similar to Strike, we use
520 Pour as a terminal skill too. In order to satisfy the pre-condition of Pour, the transform
521 between the source and target mug must belong to the family of admissible distributions.
522 We achieve the actual trajectory by designing a deterministic motion. With appropriate
523 skill level data, we can also train skill diffusion models, however, such improvement is out
524 of scope of this work.

525 **Training.** We train individual skill diffusion score-functions using the denoising score-
526 matching (DSM) loss following algorithm 2. We collect datasets of transitions observed during
527 the execution of a skill on an object and use them to train the score networks. The dataset size varies
528 according to the difficulty and diversity of a skill’s execution on a particular object. For example, we
529 need 100 successful Pick parameters for training the skill to pick the hammer and 300 successful
530 Move parameters to cover the whole workspace of the robot. For ReGrasp, we use both the Pick
531 and Move parameters.

532 **Effect of training data coverage.** If we consider “ideal” score functions and a perfect representation
533 of the factor distributions, a solution exists if there is an overlap between two connected factor
534 distributions. If such an overlapping segment does not exist, GFC will not be able to complete the
535 spatial-temporal plan. Hence, the training data for each factor (here temporal factors only) must be
536 diverse enough to ensure that the overlap exists. For example, a successful handover in *Hammer*
537 *Place* and *Hammer Strike* is not possible if the training data only consists of Pick parameters to
538 pick the hammer from the center of the handle. Similarly, if the training data for Move does not
539 cover the common workspace of both robots, our proposed algorithm will be unable to complete the
540 coordinated plan.

541 **Model architecture.** Our transformer-based score-network architecture is derived from the Dif-
542 fusion Models with Transformers (DiT) [53] implementation, also open-sourced at: [https://](https://github.com/facebookresearch/DiT)
543 github.com/facebookresearch/DiT. We follow a similar concept to that of patchifying an im-
544 age into many smaller patches, encoding each one of them using a common encoder and passing
545 it as a sequence to the transformer architecture with respective positional embeddings. In our case,
546 we consider a sequence of nodes consisting of both the object and skill parameters nodes in the
547 factor graph as the input sequence. Each node variable is encoded into a common dimension using
548 a common object node encoder and skill parameter encoder for object and skill parameter nodes
549 respectively. The output is decoded into their respective dimensions using similar decoder setup.

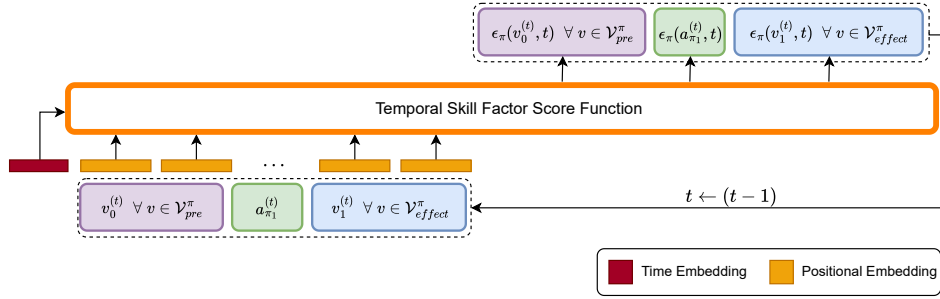


Figure S7: Transformer-based skill diffusion model. We use the noisy pre-condition, action and effect node value distribution at diffusion step t to obtain the corresponding ϵ during sampling.

Algorithm 2: Training skill score functions for a particular skill π

- 1 **Inputs:**
 - 2 Pre-condition, skill parameter and Effect nodes $(\mathcal{V}_{pre}^\pi, a_\pi, \mathcal{V}_{effect}^\pi)$
 - 3 Dataset of transitions \mathcal{D}
 - 4 Parameterized skill score function ϵ_ϕ
 - 5 Noise schedule σ
 - 6 DSM loss weight schedule λ
 - 7 **while** *not converged* **do**
 - 8 Sample batch from dataset $\mathbf{x}^{(0)} \sim \mathcal{D}$
 - 9 Sample forward diffusion timestep $t \sim [0, 1]$
 - 10 Sample Gaussian noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$
 - 11 Calculate noise coefficient σ_t
 - 12 Calculate noisy data $\mathbf{x}^{(t)} = \mathbf{x}^{(0)} + \sigma_t \epsilon$
 - 13 **end**
 - 14 Optimize parameters ϕ using:
 - 15
$$\nabla_\phi \mathbb{E}_{t, \epsilon, \mathbf{x}^{(0)}} [\lambda(t) \|\epsilon - \epsilon_\phi(\mathbf{x}^{(t)}, t)\|^2]$$
 - 16 Return $\epsilon_\pi \equiv (\text{Optimized}) \epsilon_\phi$
-

550 **Hyperparameters and computation.** We consider the hyperparameters as shown in Table S2 for
551 building our score-network.

Table S2: Hyperparameters for Score-Network with Transformer Backbone

Hyper-parameter	Value
Hidden Dimension	128
Number of Blocks	2
Number of Heads	2
MLP Ratio	2
Dropout Probability	0.1
Number of Input Channels	Varies (3-11)
Number of Output Channels	Varies (3-11)

552 For the reverse sampling steps while inference, we find the best performance using 50 steps and
 553 all results have been reported accordingly. Considering skill-object score functions with varying
 554 input nodes leads to a loss of parallel batched inference (advantage of vectorized states) and hence,
 555 an increase in computation time as compared to chaining with vectorized states. On an NVIDIA
 556 RTXTM A6000 GPU, it takes 2.6 secs for the smallest horizon task *Pour Cup* and 6 secs for the
 557 longest horizon task *Hammer Nail* to give 10 candidate node variable values. These candidates are
 558 sorted based on their extent of goal-condition satisfaction and the top 5 are selected to calculate the
 559 success performance.

560 **Example of spatial factors.** Previous work [21] considered a family of spatial factors like (left,
 561 right, top, bottom, near and far) to model collision-free object configurations. In this work,
 562 we are particularly interested in constructing a family of fixed transforms (`FixedTransform`) to
 563 model coordinated manipulation motion. For example, in order to satisfy the pre-condition of
 564 `strike(A, B)`, the transform between nodes A and B must satisfy a family of transforms sig-
 565 nifying that B must be `Aligned` with A to strike it. Thus the factor for `strike(A, B)` with
 566 `Aligned` transforms \mathcal{H}_A will look like: $f \equiv \text{distance}(\text{transform}(A, B), \mathcal{H}_A) \leq \text{permissible error}$
 567 for at least one transform. In that case, the distribution of the factor will be: $p(f = \text{True} | A, B) \propto$
 568 $\exp[-\text{distance}(\text{transform}(A, B), \mathcal{H}_A)]$. The score of such a distribution can then be calculated as

$$\epsilon_f(A^{(t)}, B^{(t)}, t) = -\nabla_{A^{(t)}, B^{(t)}} \text{distance}(\text{transform}(A^{(t)}, B^{(t)}), h_A)$$

569 where $h_A \in \mathcal{H}_A$ is the closest transform to the current transform. The distance between transforms
 570 is calculated as the summation of the Cartesian distance and the quaternion distance.

571 **S4 Real Robot Experiments**

572 **Complete setup.** We use two Franka Panda robot arms
 573 placed in parallel to demonstrate the coordinated tasks as il-
 574 lustrated in Figure S8. A pair of flexible Finray fingers [54]
 575 is attached to the parallel jaw grippers. For each of the arm,
 576 we set up a Kinect Azure camera calibrated to the origin of the
 577 arm. We use objects like mallet (hammer), stake (tent peg,
 578 nail), garden foam, a kitchen pot, two types of mugs and a
 579 rack for the considered tasks. We use segment-anything [55]
 580 and CLIP [56] to segment the objects from the RGBD image
 581 based on text descriptions and use the segmented masks to
 582 obtain the point clouds for the objects. Finally, we use ICP
 583 to align the obtained and model point clouds to calculate the
 584 transformation of the object. The procedure is done for both
 585 cameras to obtain transforms for all the detected objects in
 586 both robot’s frame of reference. For a particular object, we
 587 select the transform from the arm closest to the object to get
 588 precise pose estimation (due to better depth data). We finally
 589 use the obtained transforms to recreate the physical scene in
 590 simulation. While planning, the Frankx controller [57] is used
 to generate smooth motion toward the desired pose.

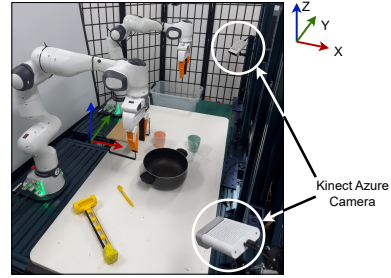


Figure S8: Real-World Experimental Setup

591 **Qualitative analysis.** We perform qualitative analysis for all four coordinated tasks using the hardware setup as shown in Figure S9, Figure S10, Figure S11 and Figure S12. We further provide detailed videos of execution in the supplementary video.

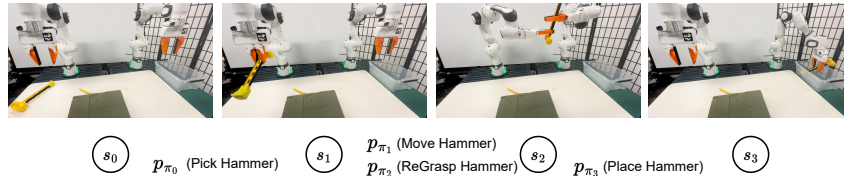


Figure S9: **Coordination task: Hammer Place** The left arm must handover the hammer to the right arm such that the hammer can be placed inside the box.

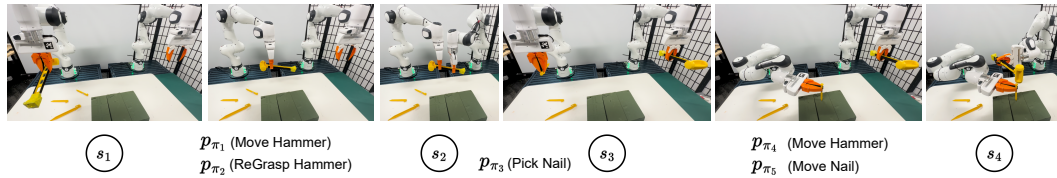


Figure S10: **Coordination task: Hammer Nail** The left arm must handover the hammer to the right arm and pick up the nail. Both arms have to coordinate in order to move the hammer and nail to a configuration in which the hammer can strike the nail.

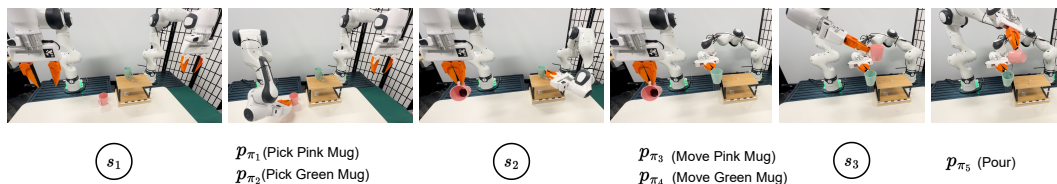


Figure S11: **Coordination task: Pour Cup** The left arm and right arm must pick up the pink mug and green mug respectively. Both arms have to coordinate in order to move the mugs to a configuration in which the left arm can pour the pink mug contents into the green mug.

594 **Failure analysis.** We try to analyze the reason for the failure of GFC in certain cases. A limiting
 595 factor of our planning framework is that the nodes denote waypoints required to be reached for

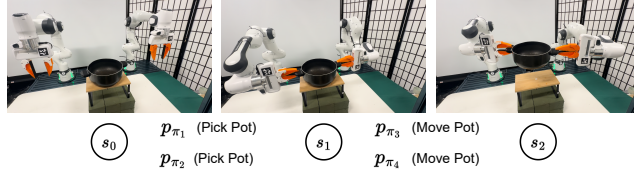


Figure S12: **Coordination task: *Bimanual Reorientation*** The left arm and right arm must pick up the pot simultaneously. Both arms have to coordinate in order to rotate the pot to a specified target reorientation angle. For the above illustration, the reorientation angle is 30deg.

596 completing the geometric execution and satisfying the goal condition without caring about the tra-
 597 jectory between them. Since we do not explicitly provide the intuition of inverse kinematics (IK) or
 598 collision, we assume that these properties are learned implicitly using the successful transitions in
 599 the training data. Hence, apart from sim-to-real gap (consisting of pose-estimation error, nature of
 600 surfaces in contact, and weight of the objects like hammer and pot), the primary reasons for failure
 601 are: (1) sampling a pose where IK cannot be computed, i.e. unreachable. (2) The sampled pose is
 602 not collision-free. We provide sim-to-real gap failures in the supplementary video.

603 **S5 More Details on Evaluation Tasks**

604 **S5.1 Hammer Nail**

605 **Task Description:** Given a scene with three boxes, a hammer in placed in one of the box covered
 606 by a lid as shown in Figure S13. There is a nail on the table. Only left arm can reach the lid, hammer
 607 and the nail. The task objective is to strike the nail by the hammer within a provided region. There
 608 is a cube in one of the boxes, picking and placing it are task-irrelevant distractions.

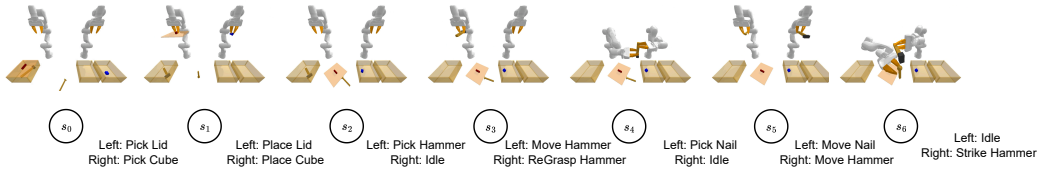


Figure S13: **Hammer Nail.** The illustration shows the *Hammer Nail* task. A successful solution to this task must complete a successful handover and coordinate to align the hammer and the nail to conduct a successful strike.

609 **What it takes to solve?** From a superficial symbolic analysis, the task can be completed if the
 610 left arm can handover the hammer to the right arm, left arm can pick up the nail to take it to the
 611 admissible region and the right arm can strike the nail by the hammer. However, the following
 612 challenges exist:

- 613 1. Hammer must be picked up and moved at a location such that the right arm can re-grasp it
 614 for a successful handover.
- 615 2. The handover must allow the right arm to satisfy the pre-condition of strike i.e. the right
 616 arm must grasp the hammer away from the head, hence the left arm must reason and pick
 617 it up by grasping close to head.
- 618 3. The re-grasp pose will affect the region where the hammer head can be reached. The
 619 left arm must reason about the hammer head’s reachability to move the nail such that the
 620 hammer and nail can be aligned.

621 **Why is this challenging?** All the above reasonings are interdependent and the effect of the initial
 622 pick pose can be seen at multiple stages of the task. This makes the task challenging as the plans
 623 fails:

- 624 1. if the initial pick pose fails to reason about handover requirements.
- 625 2. if the nail move target pose fails to satisfy the reachability of the hammer-head, which
 626 actually depends on the handover.

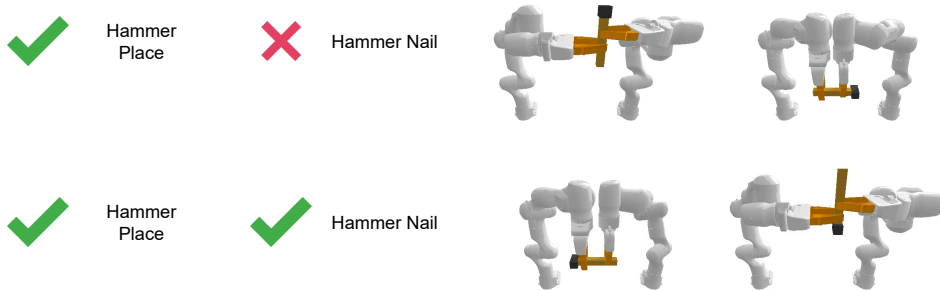


Figure S14: **Handover variations.** The hammer handover can be done in multiple ways, four of which are shown above. While placement of the hammer in the box for *Hammer Place* task can be done by re-grasping the hammer anywhere, for hammer strike in *Hammer Nail*, the hammer is encouraged to be regrasped near the tail of the handle.

627 **Failure cases:** The failures in the proposed method occur in the following situations:

- 628 1. *Method failure:* when it predicts in-feasible poses (where IK cannot be computed) or which
629 does not satisfy the pre-condition of the next skill.
- 630 2. *Trajectory planning failure:* If IK can be computed for current and target poses but no
631 collision-free trajectory can be computed (via pybullet-planning cite). This is expected as
632 GFC only solves for high-level skill transitions.
- 633 3. *Simulation failure:* While executing Pick skill, sometimes the contact vectors are noisy
634 and hence leads to pick-up failures.

635 S5.2 Bimanual Pot Reorientation

636 **Task Description:** Given a pot on a table, the task is to reorient the pot to some target orientation
637 angle (along z-axis) using two manipulators as shown in Figure S15. It is worth noting that we have
638 Pick and Move skills for individual manipulators such that we know where the pot can be grasped
639 and the reachable workspace of the manipulator.

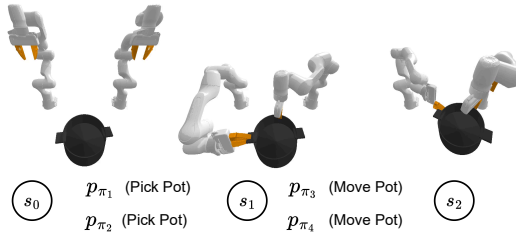


Figure S15: **Bimanual Pot Reorientation.** The task is to coordinate planning strategies to grasp a pot using two manipulators and rotate it to a target reorientation angle. The task must be done with only single-manipulator data.

640 **What it takes to solve?** This particular task can be completed if:

- 641 1. we find pick poses for both the manipulators.
- 642 2. we find feasible move poses in the workspace that satisfies the target orientation.
- 643 3. we ensure that the relative transform between two gripper poses while picking and in the
644 predicted move target poses is the same, because the grasp poses relative to the pot cannot
645 change while moving.

646 **Why is this challenging?** The task is challenging because the algorithm must decide the initial pick
647 pose by considering sequential and parallel dependencies:

- 648 1. the same pick pose relative to the pot must exist for the target reorientation angle
- 649 2. the move pose for both manipulators must satisfy both the workspace reachability for indi-
650 vidual manipulators and also have the same fixed transform as the pick poses.

651 **Failure cases:** The failures in the proposed method occur in the following situations:

- 652 1. *Method failure:* when it predicts in-feasible poses (where IK cannot be computed) or which
653 does not satisfy the fixed transform condition.
- 654 2. *Trajectory planning failure:* If IK can be computed for current and target poses but no
655 collision-free trajectory can be computed (via pybullet-planning cite). This is expected as
656 GFC only solves for high-level skill transitions.
- 657 3. *Simulation failure:* While executing Pick skill, sometimes the contact vectors are noisy
658 and hence lead to pick-up failures.

659 **S6 Extending Hammer Nail task to longer horizons**

660 In order to evaluate the extensive long-horizon planning capabilities of our proposed algorithm,
 661 we have further extended the Hammer Nail task to longer horizons as shown in Figure S16. The
 662 extended tasks particularly emphasize adding a second handover such that the hammer is handed
 663 back to the left arm after a successful hammer strike.



Figure S16: **Extension of Hammer Nail task.** We have added three new extensions to the *Hammer Nail* task. All of the new tasks focus on handling a second handover. The nature of the first handover adds further constraints into possible ways to perform the second handover. Further, we add task-irrelevant skills in between the plan skeleton to evaluate the robustness of GFC and the spatial-temporal factor graph plan representation.

664 We classify the failure cases as:

- 665 • Type 1: Method failure i.e. when the proposed algorithm fails to find suitable target param-
 666 eters.
- 667 • Type 2: Trajectory planning failure i.e. no collision-free trajectory can be computed be-
 668 tween two suitable poses.
- 669 • Type 3: Simulation failure i.e. when simulator fails to detect suitable contacts.

Table S3: Failure breakdown and task success analysis of hammer nail task and its extensions with two handovers (based on 100 trials)

Task	Task Horizon	Type 1 failure	Type 2 failure	Type 3 failure	Task Success
Hammer Nail	11	42	14	10	34
Extended Hammer Nail v1	16	43	28	5	24
Extended Hammer Nail v2	18	44	21	10	25
Extended Hammer Nail v3	20	41	25	13	21

670 Now, we show the failure breakdown and task success for all the considered *Hammer Nail* task and
671 their extensions in [Table S3](#). While we see a drop in success rates by adding a second handover to the
672 vanilla *Hammer Nail* task, GFC proved to be robust for all other task-irrelevant skills in the chain.
673 The task success of all “two handover” variants is similar even with an increasing task horizon.

674 **S7 Analyzing Inter-step dependencies**

675 Our work focuses on solving long-horizon tasks that have strong inter-step dependencies [3] and
 676 requirements for coordinated manipulation [35, 36, 37]. For example, hammering a nail not only
 677 requires extensive affordance planning to perform a handover but also requires allowing sufficient
 678 reachable workspace to align the hammer head with the nail. This also affects the success of the
 679 second handover, thus increasing the action-dependency horizon. Our framework is able to compose
 680 learned factors (diffusion models) to solve a wide variety of tasks, as long as their solutions fall in
 the combinatorial space.

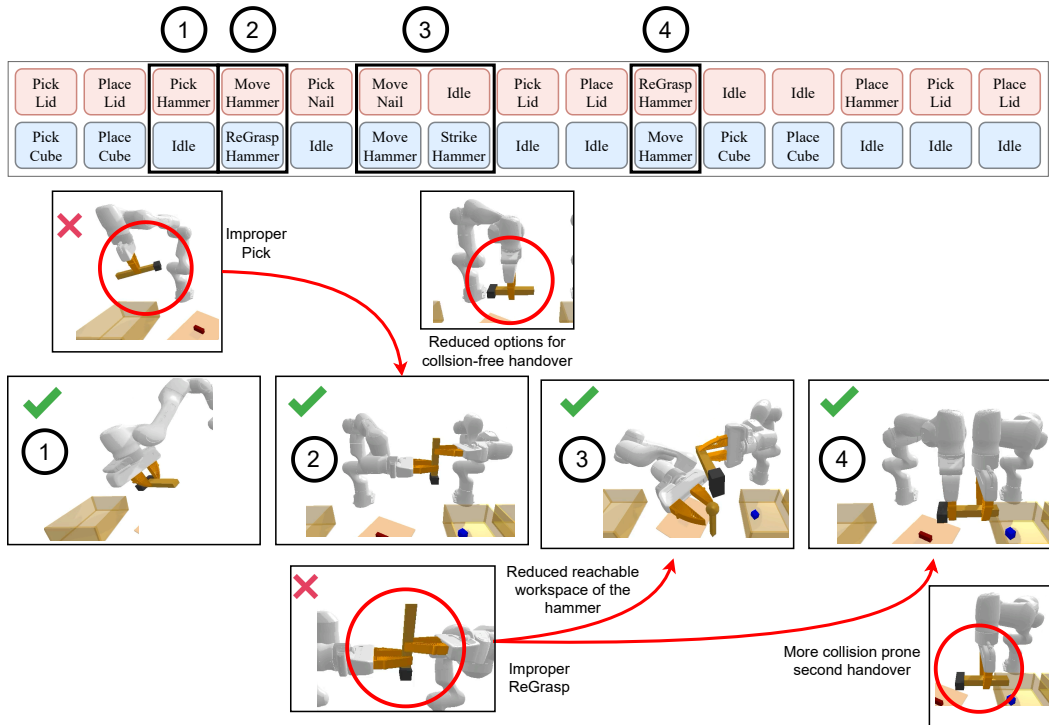


Figure S17: **Inter-step dependencies.** We show the steps and reasoning required to solve the *Hammer Nail* task. An improper initial pick can lead to a failed or unfavorable handover which might lead to difficulty in performing *Strike* and the second handover. Thus the algorithm must reason about inter-step action dependency over longer horizons to solve the task successfully.

681

682 **S8 Justifying success rates with breakdowns**

683 We elaborate on the failure and success breakdown for the vanilla *Hammer Nail* task in Table S4.
 684 Revisiting the failure categories, we classify the failure cases as:

- 685 • Type 1: Method failure i.e. when the proposed algorithm fails to find suitable target param-
 686 eters.
- 687 • Type 2: Trajectory planning failure i.e. no collision-free trajectory can be computed be-
 688 tween two suitable poses.
- 689 • Type 3: Simulation failure i.e. when simulator fails to detect suitable contacts.

Table S4: Failure breakdown and task success analysis per skill-step of hammer nail task (based on 100 trials)

Skill.No.	Skills	Type 1 failure	Type 2 failure	Type 3 failure	Accu. Success
1	Pick Lid	5	0	0	95
2	Place Lid	0	0	0	95
3	Pick Cube	0	0	0	95
4	Place Cube	6	0	0	89
5	Pick Hammer	3	0	2	84
6-7	Move Hammer - Regrasp Hammer	8	6	0	70
8	Pick Nail	4	0	8	58
9-10	Move Nail - Move Hammer	11	8	0	39
11	Hammer Strike	5	0	0	34

690 We also elaborate on the failure and success breakdown for the bimanual reorientation task in Ta-
 691 ble S5. It is worth to be noted that the skills are executed in parallel and the serialized representation
 692 of the skill sequence is shown only as a part of the analysis.

Table S5: Failure breakdown and task success analysis per skill step of bimanual pot reorientation (based on 100 trials)

Skill.No.	Skills	Type 1 failure	Type 2 failure	Type 3 failure	Accu. Success
1	Grasp Pot Left	13	0	4	83
2	Grasp Pot Right	12	0	3	68
3-4	Move Pot Left - Move Pot Right	13	12	0	53

693 We further continue the analysis for all the two handover extensions of the *Hammer Nail* task,
 694 namely for *Extended Hammer Nail v1* in Table S6, for *Extended Hammer Nail v2* in Table S7,
 695 and for *Extended Hammer Nail v3* in Table S8. We primarily note the accumulative success at the
 696 first handover, coordination for the hammer Strike, and the second handover. With an increasing
 697 task horizon, the proposed approach is invariant to task-irrelevant distractions and maintains similar
 698 success.

Table S6: Failure breakdown and task success analysis per skill-step of hammer nail task extension v1 with two handovers (based on 100 trials)

Skill.No.	Skills	Type 1 failure	Type 2 failure	Type 3 failure	Accu. Success
1	Pick Lid	4	0	0	96
2	Place Lid	0	0	0	96
3	Pick Cube	0	0	0	96
4	Place Cube	5	0	0	91
5	Pick Hammer	4	0	2	85
6-7	Move Hammer - Regrasp Hammer	11	13	0	61
8	Pick Nail	3	0	3	55
9-10	Move Nail - Move Hammer	7	9	0	39
11	Hammer Strike	3	0	0	36
12-13	Move Hammer - Regrasp Hammer	4	6	0	26
14	Place Hammer	0	0	0	26
15	Pick Lid	2	0	0	24
16	Place Lid	0	0	0	24

Table S7: Failure breakdown and task success analysis per skill-step of hammer nail task extension v2 with two handovers and some task-irrelevant skills (based on 100 trials)

Skill No.	Skills	Type 1 failure	Type 2 failure	Type 3 failure	Accu. Success
1	Pick Lid	4	0	0	96
2	Place Lid	0	0	0	96
3	Pick cube	0	0	0	96
4	Place Cube	4	0	0	92
5	Pick Hammer	5	0	2	85
6-7	Move Hammer - Regrasp Hammer	12	14	0	59
8	Pick Nail	2	0	1	56
9-10	Move Nail - Move Hammer	4	0	7	45
11	Hammer Strike	1	0	0	44
12	Pick Lid	3	0	0	41
13	Place Lid	0	0	0	41
14-15	Move Hammer - Regrasp Hammer	6	7	0	28
16	Place Hammer	0	0	0	28
17	Pick Lid	3	0	0	25
18	Place Lid	0	0	0	25

Table S8: Failure breakdown and task success analysis per skill-step of hammer nail task extension v3 with two handovers and many task-irrelevant skills (based on 100 trials)

Skill No.	Skills	Type 1 failure	Type 2 failure	Type 3 failure	Accu. Success
1	Pick Lid	5	0	0	95
2	Place Lid	0	0	0	95
3	Pick cube	0	0	2	93
4	Place Cube	4	0	0	89
5	Pick Hammer	3	0	2	84
6-7	Move Hammer - Regrasp Hammer	4	8	0	72
8	Pick Nail	3	0	6	63
9-10	Move Nail - Move Hammer	7	9	0	47
11	Hammer Strike	5	0	0	42
12	Pick Lid	1	0	2	39
13	Place Lid	0	0	0	39
14-15	Move Hammer - Regrasp Hammer	5	8	0	26
16	Pick cube	0	0	0	26
17	Place Cube	1	0	0	25
18	Place Hammer	3	0	0	22
19	Pick Lid	0	0	1	21
20	Place Lid	0	0	0	21