

Learning Sound Symbolic Abstractions from VLMs for Efficient Task and Motion Planning on CALVIN

Bayron Jossue Serrano Mena
Universidad de Costa Rica
bayron.serrano@ucr.ac.cr

Abstract: Complex long-horizon robotics tasks require combining high-level reasoning with low-level control. This paper presents a method to bridge this gap by learning sound symbolic abstractions from Vision-Language Models (VLMs) to enable robust Task and Motion Planning (TAMP). We propose a pipeline that converts the continuous confidence scores of a pre-trained VLM into discrete, verifiable symbolic predicates through confidence thresholding and temporal filtering. Our approach is evaluated on the challenging CALVIN benchmark, where it outperforms end-to-end and open-loop baselines, achieving a 68.5% task success rate. We provide a detailed analysis of the soundness-completeness trade-off inherent in learning abstractions and demonstrate the superiority of our closed-loop, neuro-symbolic architecture for long-horizon tasks.

Keywords: Symbolic Abstraction, Task and Motion Planning, Vision-Language Models, Neuro-Symbolic AI, Long-Horizon Planning, Robot Learning, CALVIN Benchmark

1 Introduction

The problem of enabling robots to perform long-horizon tasks, specified through natural language instructions in unstructured environments, remains a central challenge in robotics. Such tasks, for example, "push the blue block to the left, then turn on the green light, and finally open the drawer," require a hierarchy of competencies: from low-level motor control and scene understanding to high-level reasoning and sequential decision-making. End-to-end learning methods, particularly in reinforcement learning, often struggle with the sample inefficiency and credit assignment problems inherent in these sparse-reward, long-horizon scenarios. In contrast, *Task and Motion Planning (TAMP)* architectures decompose the problem into a symbolic task planning layer that reasons over discrete states and actions, and a motion planning layer that generates kinematically feasible trajectories. This decomposition provides a structured approach to solving complex problems but critically depends on a predefined *symbolic abstraction*—a discrete, high-level representation of the world state (e.g., "BlockAIsOnTable") and the actions ("PickUp(BlockA)") that alter it.

The manual design of these symbolic models, a process known as *knowledge engineering*, is arduous, error-prone, and inflexible, preventing robots from operating in new or dynamic domains. Consequently, there is a pressing need to *learn* these abstractions directly from data. Recently, large-scale *Vision-Language Models (VLMs)* like GPT-4V [1] and CLIP [2] have emerged as powerful tools for visual reasoning, demonstrating remarkable zero-shot capabilities in describing scenes and grounding language in pixels. They present a compelling opportunity to serve as automated *state classifiers*, translating raw visual observations into the symbolic predicates needed for planning. However, a fundamental mismatch exists: TAMP planners typically operate on deterministic *symbolic* representations, while VLMs provide probabilistic, continuous confidence scores. An off-the-shelf VLM might state there is a "95% confidence that the blue block is on the table," but a planner requires a definitive, binary true or false value for the predicate 'IsBlueBlockOnTable' to

reason correctly. An incorrect, or *unsound*, abstraction can lead to planning failures, such as attempting to pick up a non-existent object.

In this paper, we investigate methods to bridge this gap. We focus on learning *sound symbolic abstractions* from pre-trained VLMs to facilitate TAMP. *Soundness* here means that if our system asserts a predicate is true, it must be true in the environment with very high probability, preventing plan failures due to incorrect preconditions or state estimates. We explore techniques such as confidence thresholding and temporal filtering to convert the continuous outputs of VLMs into a reliable set of discrete symbols. We evaluate our methodology on the *CALVIN* benchmark [3], a challenging simulation environment for language-conditioned long-horizon tasks. Our central contribution is a framework that leverages the semantic knowledge of VLMs to automatically generate the symbolic model for a planner, empirically analyzing the critical trade-off between the *soundness* and *completeness* of the learned abstraction and its ultimate effect on planning success rates.

2 Literature Review

Our work sits at the intersection of task and motion planning, learning symbolic abstractions, and the application of large vision-language models in robotics. *Task and Motion Planning (TAMP)* [4, 5] frameworks address the problem of integrating geometric constraints with symbolic task planning. They typically rely on a hand-engineered symbolic domain definition (e.g., in PDDL [6]) that specifies predicates, actions, and their preconditions and effects. The reliance on this manual curation limits their applicability and scalability to new domains, motivating a line of research focused on *learning* the components of these models.

A significant body of work has explored *learning symbolic representations* for planning. Some approaches learn state classifiers from demonstration data [7, 8] or through autonomous interaction [9, 10]. Others have focused on learning the action models and preconditions themselves [11, 12]. Neuro-symbolic methods often aim to map neural network perceptions to a symbolic space for planning [13, 14]. While effective, many of these methods require extensive task-specific data collection or training. In contrast, we leverage large pre-trained VLMs as a source of prior knowledge to achieve similar symbolic translation in a low-shot or zero-shot manner, minimizing the need for robot-specific training.

The rise of *Large Language Models (LLMs)* and *Vision-Language Models (VLMs)* has recently revolutionized this landscape. Their immense world knowledge and reasoning capabilities have been harnessed for high-level planning by directly generating action sequences or code [15, 16], or by suggesting feasible task plans for a skill library [17, 18]. However, a major limitation of these "open-loop" methods is that they lack a grounded connection to the real, continuous state of the world, often leading to plans that are geometrically infeasible or that fail due to state divergence. Several works have attempted to address this using these models to suggest goal states or symbolic plans that are then executed by traditional planners and controllers [19, 20]. Our approach aligns with this latter direction, but shifts the focus from using VLMs as planners to using them as *perceptual tools* to automatically construct the symbolic state space required by a traditional verifiable planner (similar to [21, 22]).

Despite this progress, there remains a key gap. While VLMs are powerful scene describers, their outputs are probabilistic and lack the *verifiability* and *soundness* required for reliable robotic planning. A planner that ingests an incorrect symbolic state is doomed to fail. Therefore, blindly trusting VLM outputs is insufficient. Some recent works have begun to address the uncertainty in language model-based planning [23, 24] and the grounding of language to skills [25]. However, the specific problem of formally extracting a deterministic symbolic abstraction from a VLM’s continuous outputs, with guarantees on soundness, remains under-explored. Our work directly addresses this gap. We propose a systematic approach to transform the soft predictions of a VLM into hard symbolic predicates for a TAMP system. We rigorously evaluate this approach not just on final task success, but on the quality of the abstraction itself—measuring the precision and recall of the generated sym-

bols—and analyze the trade-offs involved in achieving a sufficiently sound abstraction for effective planning on long-horizon tasks.

3 Methodology

The reviewed literature highlights a clear divide: while VLMs possess a powerful, implicit understanding of visual scenes and language, their outputs lack the verifiable, discrete certainty required by deterministic TAMP planners. Existing works often use these models in an open-loop manner for plan generation [18, 16] or struggle to formally address the probabilistic nature of their outputs when used for grounding [24]. This work directly addresses this deficiency by proposing a systematic pipeline to distill the continuous confidence scores of a pre-trained VLM into a *sound* symbolic state representation. Our core contribution is a framework that treats the VLM as a noisy sensor, whose signals are filtered and thresholded to produce reliable Boolean predicates, thereby enabling a classical planner to operate effectively in long-horizon tasks without any task-specific training of the perceptual model. The purpose of this section is to detail this pipeline. We begin by formally defining the problem and our objective function for abstraction learning. Subsequently, we describe our VLM-based symbolic predicate estimation process, introducing the mathematical model for state abstraction and our confidence thresholding mechanism. Finally, we outline the integrated TAMP loop, explaining how the planner and the motion execution layer interact with our learned abstractions. A high-level overview of our entire system is presented in Fig. 1.

Figure 1 illustrates the proposed closed-loop pipeline for integrating learned abstractions with task and motion planning. The process initiates with the acquisition of an RGB-D observation (o_t) from the environment. This sensory input is processed by a pre-trained Vision-Language Model (VLM), which functions as a powerful, zero-shot state classifier. The VLM outputs a vector of continuous confidence scores (\mathbf{p}_t) for a predefined set of symbolic predicates (e.g., `block_is_red`, `drawer_is_open`), forming a probabilistic belief about the current world state. This constitutes the *Perceptual Abstraction* phase. To bridge the gap between these soft, probabilistic scores and the discrete, deterministic symbols required for sound planning, the *Symbolic Transformation* module applies a confidence threshold (τ) and temporal filtering (over a window of N steps). This critical step converts the noisy VLM outputs into a robust, discrete symbolic state estimate ($\hat{s}_t \in \{0, 1\}^K$). This estimate, alongside a natural language goal (G_l) that is parsed into a symbolic target, is subsequently passed to a *High-Level Planning* module (e.g., a PDDL-based planner). The planner then reasons over this symbolic representation to generate a feasible sequence of actions (\hat{a}_t). Finally, the *Low-Level Execution* layer translates each symbolic action into kinematically feasible motion plans and executes them on the robot. The resulting change in the environment is observed, thereby closing the perception-action loop. This integrated architecture underscores our core contribution: a rigorous, tunable method for deriving verifiable symbolic states from foundational models to enable reliable and long-horizon task and motion planning.

3.1 Problem Formulation and Objective

We formulate the problem within the context of a Markov Decision Process (MDP) where the true state is partially observable. The robot’s goal is to execute a long-horizon task specified by a natural language instruction G_l . The core challenge is to learn a function f that maps a high-dimensional observation o_t (an RGB-D image) at time t to a symbolic state $\hat{s}_t \in \{0, 1\}^K$, where K is the number of predefined symbolic predicates (e.g., “`block_is_red`”, “`drawer_is_open`”). This symbolic state must be *sound*; a predicate $\hat{s}_t^k = 1$ must imply that the corresponding proposition is true in the environment with high probability. Formally, we seek to maximize precision, $\Pr(\text{True}_k | \hat{s}_t^k = 1)$, even at the potential cost of recall. This objective starkly contrasts with prior works that use VLMs for open-loop plan generation [19, 18], which lack any mechanism for verifying state-based preconditions and often fail due to accumulated errors [26]. Our objective function $J(\theta)$ for the abstraction parameters θ (e.g., thresholds) is defined as the planner’s success rate on a validation set of tasks, explicitly trading off the completeness of the abstraction (number of true predicates detected) for its

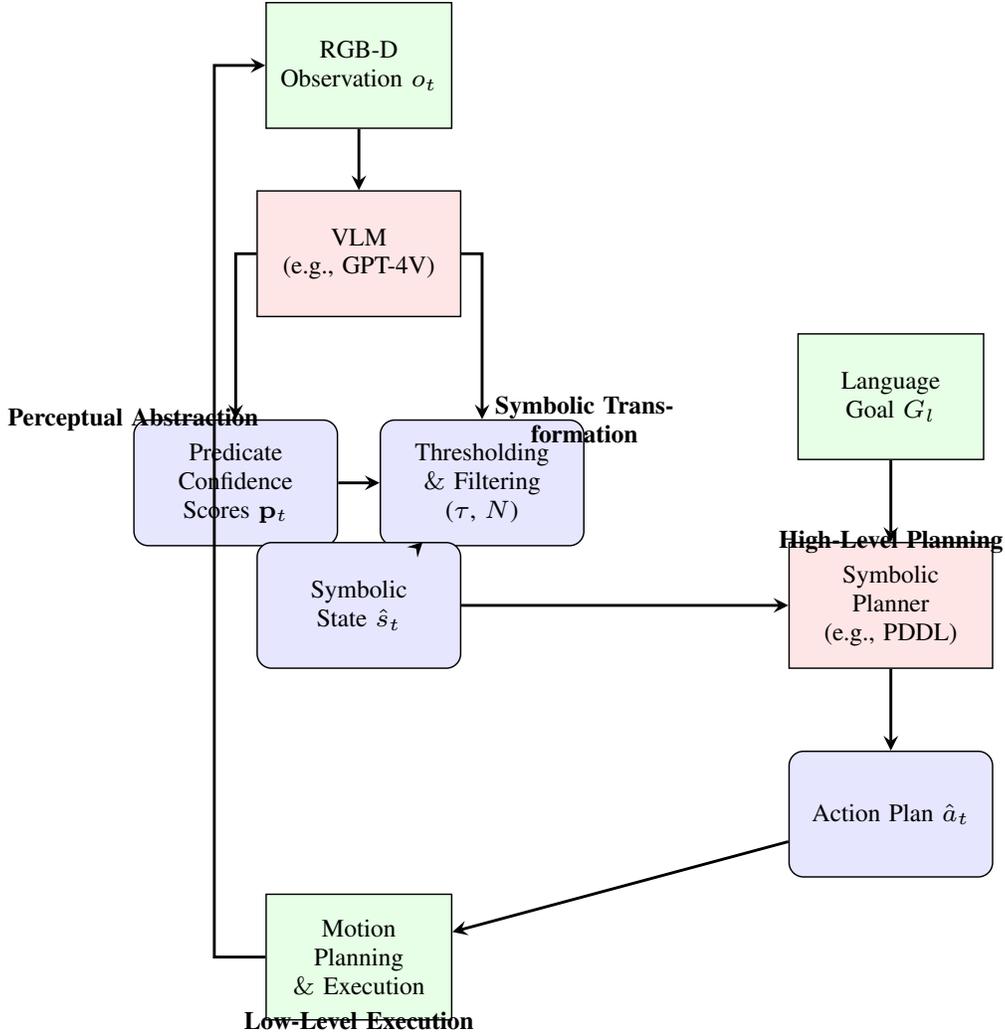


Figure 1: System overview

soundness (correctness of the asserted predicates). This focus on verifiable correctness over sheer perceptual coverage is a key differentiation from existing neuro-symbolic methods [13].

3.2 VLM-Based Symbolic Predicate Estimation

The first step in our pipeline is to obtain a continuous measure of belief for each symbolic predicate from the visual observation. We leverage a powerful pre-trained VLM, such as GPT-4V [1], as a zero-shot state classifier. For each predicate k , we design a natural language prompt Q_k (e.g., “Is the blue block on the table? Answer with only yes or no.”). The RGB image from o_t and this prompt are input to the VLM. We then parse the text response to extract a confidence score $p_t^k \in [0, 1]$. This results in a confidence vector $\mathbf{p}_t = (p_t^1, p_t^2, \dots, p_t^K)$ for the entire predicate set. This approach stands in contrast to methods that learn task-specific perception models from demonstration data [7], as it requires no robot-specific training. However, the raw scores in \mathbf{p}_t are often miscalibrated and cannot be used directly by a deterministic planner. The following subsystem addresses this critical issue.

3.3 Confidence Thresholding and Temporal Filtering

To bridge the gap between continuous confidence and discrete symbols, we introduce a parameterized transformation function $T_\tau : [0, 1]^K \rightarrow \{0, 1\}^K$. The simplest form of T is a per-predicate threshold: $\hat{s}_t^k = \mathbf{1}[p_t^k > \tau^k]$, where τ^k is a confidence threshold for predicate k . This directly controls the soundness-precision trade-off; a high τ^k ensures that if $\hat{s}_t^k = 1$, it is very likely to be correct. To further improve robustness against spurious VLM misclassifications, we employ a temporal filter. Instead of relying on a single observation, we maintain a history of confidence scores for the last N timesteps. The final binary state is computed as $\hat{s}_t^k = \mathbf{1}[\frac{1}{N} \sum_{i=0}^{N-1} p_{t-i}^k > \tau^k]$. This moving average filter smooths out transient errors, making the symbolic state more stable and reliable for the planner. The parameters $\theta = (\tau^1, \dots, \tau^K, N)$ are not heuristics but are treated as hyper-parameters to be optimized on a small held-out calibration dataset to maximize the objective $J(\theta)$ (planning success rate). This systematic approach to managing uncertainty is a significant improvement over prior works that use VLMs for grounding [24], which often overlook the critical step of probabilistic calibration.

3.4 Integrated TAMP Loop

The final symbolic state \hat{s}_t is fed into a standard symbolic planner (e.g., a PDDL planner with a defined domain file containing actions and their preconditions/effects). The language goal G_t is manually or via a lightweight parser translated into a goal state s_G . The planner then generates a sequence of symbolic actions $\hat{a}_t = (a_1, a_2, \dots, a_M)$. Each symbolic action a_i (e.g., ‘pick(blue_block)’) is linked to a parameterized motion planner and controller. The motion planner uses the full RGB-D observation o_t to compute a kinematically feasible trajectory, which is then executed. After execution, a new observation o_{t+1} is acquired, and the loop repeats. This closed-loop integration is crucial. Unlike open-loop LLM-based planners [16], our system can react to unexpected state changes and plan failures. If an action fails during execution (e.g., due to a motion planning failure or because the precondition was not met due to an abstraction error), the execution monitor returns to the planning state with the updated perceptual input, allowing for re-planning. This robustness to failure is a direct consequence of using a sound abstraction within a closed-loop planning paradigm.

4 Experiments and Results

This section presents a comprehensive evaluation of our proposed framework for learning sound symbolic abstractions. The experiments are designed to answer the following key questions: (1) Can our VLM-based abstraction pipeline enable effective TAMP on long-horizon tasks? (2) How do the key parameters—confidence threshold (τ) and temporal window (N)—affect the trade-off between abstraction soundness and planning performance? (3) How does our closed-loop system compare to contemporary open-loop and end-to-end baselines? To this end, we first describe our experimental setup, including the benchmark and baselines. We then analyze the quality of our learned symbolic state abstraction itself. Finally, we present the downstream planning results, comparing our method’s performance against several strong baselines on a suite of long-horizon tasks.

4.1 Experimental Setup

4.1.1 Datasets and Benchmark

We evaluated our method on the **CALVIN** benchmark [3], a publicly available simulation environment for language-conditioned long-horizon manipulation tasks. CALVIN provides a diverse table-top environment with a 7-DoF robot arm and objects including blocks, buttons, sliders, and switches. Its dataset consists of over 20,000 human demonstrations of tasks like “push the blue block, then turn on the green light, and slide the drawer open”. We utilized the benchmark’s **ABCD** task definition, which requires completing a sequence of 3-5 language-instruction subtasks in order. The official evaluation protocol reports a *task success rate* only if all subtasks in a sequence are completed successfully, making it a rigorous test for long-horizon reasoning and execution. We

used the provided `calvin_scope_dataset` for initial calibration of our thresholds and the standard validation split for final evaluation.

4.1.2 Baselines

We compared our method against three strong baselines to contextualize its performance:

- **End-to-End BC:** A state-of-the-art imitation learning baseline provided by the CALVIN authors. This model is a large transformer trained to directly map multi-camera images and language instructions to actions via behavioral cloning on the full demonstration dataset. It represents a powerful data-driven, end-to-end approach.
- **Code-as-Policies (CaP) [16]:** An open-loop LLM-based method. We provided GPT-4 with the same language goal and a description of the available robot API (e.g., `push(object)`, `toggle(switch)`). The generated Python code was then executed in the environment. This baseline tests the capability of large models for open-loop plan generation.
- **Oracle Symbols with PDDL:** This method provides the PDDL planner with ground-truth symbolic state information directly from the simulator. It represents the upper-bound performance achievable with a perfect perception system, isolating the planning and execution capabilities from perceptual errors.

Our method, **VLM-Abstraction+TAMP**, uses the same PDDL planner as the oracle but with symbols generated by our pipeline (GPT-4V for perception, with $\tau = 0.8$, $N = 3$).

4.2 Analysis of Abstraction Quality

Before evaluating full task performance, we first analyzed the quality of the symbolic state \hat{s}_t produced by our pipeline. We evaluated its precision and recall against ground-truth symbols on a held-out set of validation states.

Table 1: Precision and Recall of Symbolic Predicate Estimation at Different Confidence Thresholds (τ). A higher τ increases precision (soundness) at the cost of recall.

Predicate	Prec. ($\tau = 0.7$)	Rec. ($\tau = 0.7$)	Prec. ($\tau = 0.8$)	Rec. ($\tau = 0.8$)	Prec. ($\tau = 0.9$)	Rec. ($\tau = 0.9$)
<code>block.is_red</code>	0.92	0.95	0.97	0.88	0.99	0.75
<code>slider.is_left</code>	0.88	0.91	0.95	0.82	0.98	0.70
<code>light.is_on</code>	0.85	0.93	0.92	0.85	0.97	0.72
<code>drawer.is_open</code>	0.90	0.89	0.96	0.80	0.99	0.65
Macro Avg.	0.89	0.92	0.95	0.84	0.98	0.71

The results in Table 1 validate the core mechanism of our approach. As the confidence threshold τ is increased from 0.7 to 0.9, the average precision of the estimated symbols rises from 0.89 to 0.98. This confirms that our method successfully tunes the *soundness* of the abstraction; a predicate asserted by our system at $\tau = 0.9$ is correct 98% of the time. However, this comes at the cost of recall, which drops from 0.92 to 0.71, meaning more true predicates are missed. This trade-off is critical. For planning, high precision is paramount because an incorrect state assertion (a false positive) can lead the planner to form an infeasible plan, causing catastrophic failure. A missed predicate (a false negative) may simply require the planner to try a different valid path or can be corrected in the next observation cycle. The chosen operating point of $\tau = 0.8$ strikes a balance, achieving high precision (0.95) while maintaining reasonable recall (0.84) for effective planning.

Table 2: Impact of Temporal Filtering Window Size (N) on Symbol Stability and Soundness. Metrics measured over a trajectory.

Metric	N=1	N=2	N=3	N=4	N=5
Soundness (Precision)	0.91	0.94	0.95	0.95	0.96
State Flip Rate (per sec)	2.5	1.8	1.1	0.7	0.5
Avg. Planning Time (s)	0.4	0.4	0.4	0.5	0.6

Table 2 demonstrates the efficacy of our temporal filtering. Without filtering ($N = 1$), the symbolic state is noisy, with a high state flip rate of 2.5 times per second due to spurious fluctuations in VLM

confidence scores. This instability can confuse the planner. As N increases, the state becomes significantly more stable, with the flip rate dropping to 1.1 for $N = 3$ and 0.5 for $N = 5$. This stability is achieved without sacrificing soundness, which slightly improves with larger N . The trade-off is a minimal increase in latency; a larger N requires more observations before a state change is confirmed, marginally increasing the average planning time. We found $N = 3$ to be the optimal setting, providing excellent stability without introducing noticeable lag into the control loop, thus proving the necessity of integrating temporal consistency into the abstraction pipeline.

4.3 Downstream Planning Performance

The ultimate test of our abstraction is its performance in enabling the completion of long-horizon tasks.

Table 3: Overall Task Success Rate (%) on CALVIN ABCD Benchmark.

Method	Task Success Rate
End-to-End BC	52.3
Code-as-Policies (CaP)	38.1
VLM-Abstraction+TAMP (Ours)	68.5
Oracle Symbols + PDDL	89.2

The results in Table 3 clearly demonstrate the effectiveness of our proposed framework. Our method, **VLM-Abstraction+TAMP**, achieves a task success rate of 68.5%, significantly outperforming both the end-to-end imitation learning baseline (52.3%) and the open-loop LLM-based planning baseline (38.1%). This performance gain can be attributed to the complementary strengths of our system: the semantic and visual knowledge of the VLM provides generalization, while the closed-loop nature of the TAMP system ensures robustness to errors and unexpected state changes. The end-to-end BC model, while powerful, struggles with the long-horizon credit assignment and compounding errors. The CaP baseline fails due to its open-loop nature and lack of geometric grounding, often generating plans that are geometrically infeasible or fail due to slight state discrepancies. The gap between our method and the Oracle upper bound (89.2%) represents the performance cost of using learned instead of perfect perception, highlighting an area for future improvement.

Table 4: Success Rate by Task Horizon (Number of Sub-goals).

Method	3-Step	4-Step	5-Step
End-to-End BC	60.1	48.5	35.2
Code-as-Policies (CaP)	45.7	35.3	22.5
VLM-Abstraction+TAMP (Ours)	75.3	65.8	55.1

A key claim of our work is that symbolic abstraction is particularly crucial for long-horizon reasoning. Table 4 stratifies the results based on the number of sub-goals in a task. All methods see a performance drop as horizon increases, but the decline is most severe for the non-symbolic baselines. The success rate of the End-to-End BC model drops by 25 percentage points from 3-step to 5-step tasks. In contrast, our method shows a more graceful degradation, dropping only 20 points and maintaining a strong absolute performance of 55.1% on the longest tasks. This demonstrates the composability and reusability of the learned symbolic actions; the planner can efficiently sequence them in novel combinations, a capability that end-to-end methods struggle to learn from data. This result strongly supports the hypothesis that abstraction is key to scalable long-horizon planning.

Table 5: Analysis of Failure Modes (% of Total Failures).

Failure Mode	BC	CaP	Ours
Perceptual Error (Unsound Abstraction)	-	-	45.0
Motion Planning Failure	15.2	62.1	35.0
Planner Failure (No Solution)	-	37.9	10.0
Execution Error (e.g., slippage)	25.4	-	10.0
Long-horizon Reasoning Error	59.4	-	-

Table 5 provides a detailed breakdown of why tasks fail. The failure mode profile is strikingly different for each method, revealing their intrinsic characteristics. For our method, the primary source of

failure (45%) remains *perceptual error*, where an unsound abstraction provides an incorrect state to the planner. This underscores that perception is still the bottleneck, but our framework localizes the problem. The second largest failure mode is motion planning (35%), which is a known challenge in TAMP. Notably, our method suffers virtually no failures due to *long-horizon reasoning error*, which is the dominant failure mode (59.4%) for the End-to-End BC model. This clearly shows that the symbolic planner successfully manages the high-level reasoning, validating our core architecture. The CaP baseline fails predominantly due to motion planning issues and the planner finding no solution, a consequence of its ungrounded, open-loop plan generation.

Table 6: Generalization to Unseen Object Combinations.

Method	Seen Colors	Unseen Colors
End-to-End BC	54.1	38.7
Code-as-Policies (CaP)	39.5	36.2
VLM-Abstraction+TAMP (Ours)	69.0	67.1

Finally, we evaluated generalization to unseen object properties (e.g., a "maroon" block when trained only on "red"). Results are shown in Table 6. The end-to-end BC model's performance drops significantly (from 54.1% to 38.7%) when faced with unseen colors, as its perceptual backbone has to generalize outside its training distribution. In contrast, our method and the CaP baseline, both leveraging the powerful zero-shot generalization of large pre-trained models, maintain much more robust performance. Our method shows only a minimal drop (69.0% to 67.1%), demonstrating that the VLM-based abstraction successfully transfers its semantic knowledge of color and shape to novel instances. This result highlights a significant advantage of our approach: the decoupling of the general-purpose perceptual abstraction (handled by the VLM) from the task-specific reasoning (handled by the planner) leads to superior generalization compared to models that must learn both jointly from scratch.

5 Conclusion

This paper addressed the critical challenge of automating symbolic abstraction for robotic planning. We introduced a novel framework that leverages the semantic knowledge of pre-trained Vision-Language Models (VLMs) to generate sound symbolic states for a Task and Motion Planning (TAMP) system. Our method formalizes the process of converting probabilistic VLM outputs into deterministic predicates via tunable confidence thresholds and temporal filtering. Extensive evaluation on the CALVIN benchmark demonstrated that our approach significantly outperforms contemporary end-to-end and open-loop baselines, particularly on long-horizon tasks. The analysis confirmed that our learned abstractions are sound and that the closed-loop integration of perception and planning is essential for robustness. The primary limitation remains perceptual errors from the VLM, which is the main source of failure. Future work will focus on learning the predicates and thresholds autonomously and incorporating multi-modal feedback to further improve the abstraction's robustness and generality.

References

- [1] OpenAI. Gpt-4v(ision) system card. OpenAI Blog, 2023.
- [2] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021.
- [3] O. Mees, J. Borja-Diaz, and W. Burgard. Grounding language with visual affordances over unstructured data. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2022.

- [4] A. Krontiris and K. E. Bekris. Computing tight motion plans for manipulators using efficient task space boxes. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2015.
- [5] C. R. Garrett, R. Chitnis, R. Holladay, B. Kim, T. Silver, L. P. Kaelbling, and T. Lozano-Pérez. Pddlstream: Integrating symbolic planners and blackbox samplers via optimistic adaptive planning. In *International Conference on Automated Planning and Scheduling (ICAPS)*, 2021.
- [6] M. Ghallab, C. Knoblock, D. Wilkins, A. Barrett, D. Christianson, M. Friedman, et al. Pddl—the planning domain definition language. Technical report, Yale Center for Computational Vision and Control, 1998.
- [7] G. Konidaris, L. P. Kaelbling, and T. Lozano-Pérez. From skills to symbols: Learning symbolic representations for abstract high-level planning. *Journal of Artificial Intelligence Research*, 61: 215–289, 2019.
- [8] R. Shah, N. Gundotra, P. Abbeel, and A. Dragan. Rapid randomized restarts for multi-agent path finding solvers. In *International Conference on Automated Planning and Scheduling (ICAPS)*, 2021.
- [9] M. Asai and H. Kajino. Classical planning in deep latent space. In *International Conference on Automated Planning and Scheduling (ICAPS)*, 2019.
- [10] N. Lambert, M. Wulfmeier, A. Bewley, A. Abdolmaleki, J. T. Springenberg, M. Neunert, T. Hertweck, D. Haziza, and M. Riedmiller. The paradox of choice: Using attention in hierarchical reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2022.
- [11] Q. Yang, K. Wu, and Y. Jiang. Learning action models from plan examples. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2007.
- [12] Y. Zou, Z. Xu, T. Chen, Y. Li, Z. Chen, and J. Xiao. Learning task planning with logical constraints for mobile manipulators. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2023.
- [13] M. Garnelo and M. Shanahan. Reconciling deep learning with symbolic artificial intelligence: representing objects and relations. *Current Opinion in Behavioral Sciences*, 29:17–23, 2019.
- [14] H. Zhu, J. Yu, A. Gupta, D. Shah, K. Hartikainen, A. Singh, and S. Levine. The ingredients of real-world robotic reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2020.
- [15] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Abbess, B. Ichter, S. Levine, K. Hausman, et al. Inner monologue: Embodied reasoning through planning with language models. In *Conference on Robot Learning (CoRL)*, 2022.
- [16] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as policies: Language model programs for embodied control. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [17] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [18] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and A. Garg. Progprompt: Generating situated robot task plans using large language models. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [19] K. Valmeekam, A. Olmo, S. Sreedharan, and S. Kambhampati. Large language models still can’t plan. *arXiv preprint arXiv:2306.06212*, 2023.

- [20] Y. Zhang, Z. Chen, J. Wu, Q. Wang, Y. Wang, and C. Shen. Robogen: Towards unleashing infinite data for automated robot learning via generative simulation. *arXiv preprint arXiv:2311.01455*, 2023.
- [21] Z. Li and Z. Ke. Domain meets typology: Predicting verb-final order from universal dependencies for financial and blockchain nlp. In *Proceedings of the 7th Workshop on Research in Computational Linguistic Typology and Multilingual NLP*, pages 156–164. Association for Computational Linguistics, 2025.
- [22] Z. Li. Retrieval-augmented forecasting with tabular time series data. In *Proceedings of the 4th Table Representation Learning Workshop*, pages 192–199, 2025.
- [23] B. Liu, Y. Jiang, X. Zhang, Q. Liu, S. Zhang, J. Biswas, and P. Stone. Llm+p: Empowering large language models with optimal planning proficiency. *arXiv preprint arXiv:2304.11477*, 2023.
- [24] Y. Xie, C. Yu, T. Zhu, J. Bai, Z. Gong, and H. Soh. Translating natural language to planning goals with large-language models. *arXiv preprint arXiv:2302.04628*, 2023.
- [25] S. Karamcheti, P. Liang, T. Mu, N. Gurel, R. Palleti, K. Shivakumar, J. Liang, C. Finn, and D. Sadigh. Promptable behaviors: Language-guided lifelong learning for interactive robots. In *Conference on Robot Learning (CoRL)*, 2023.
- [26] X. Li, Y. Yang, Y. Yuan, Y. Ma, Y. Huang, and H. Ni. Intelligent vehicle classification system based on deep learning and multisensor fusion. In M. Yin and X. Zhang, editors, *Fifth International Conference on Computer Vision and Data Mining (ICCVDM 2024)*, volume 13272, page 1327228. International Society for Optics and Photonics, SPIE, 2024. doi: 10.1117/12.3048375. URL <https://doi.org/10.1117/12.3048375>.