# Discovering High Level Patterns from Simulation Traces

SEAN MEMERY, University of Edinburgh, United Kingdom

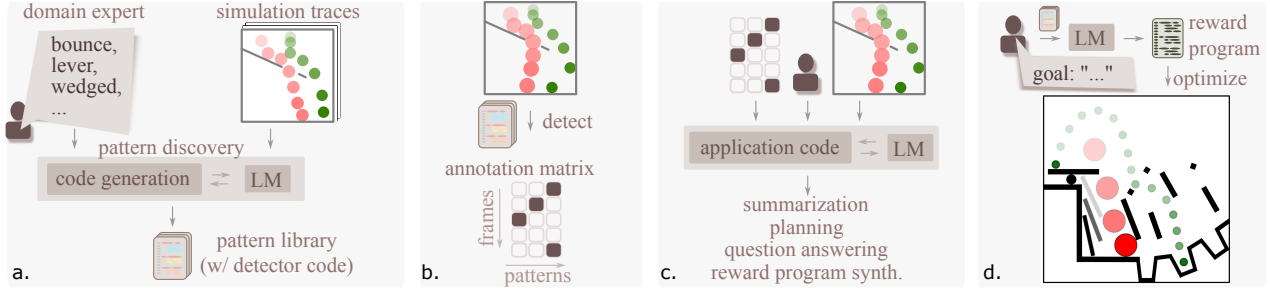KARTIC SUBR, University of Edinburgh, United Kingdom

Fig. 1. Our method *discovers high-level patterns from low-level simulation traces*, improving LM physical reasoning without fine-tuning. (a) Pattern discovery utilises evolutionary programming to detect coarse-grained patterns (e.g., 'a lever launches a ball') from raw simulation states. (b) A library of detecting code can annotate a simulation trace, creating a matrix of pattern activations (c) Annotated traces are useful for downstream tasks: *summarization, physics planning, question answering, and reward program synthesis.* (d) An example reward program from a natural language goals "*Get the green ball in the second bucket, by passing between the two obstacles*", optimized to solve the task.

Artificial intelligence (AI) agents embedded in environments with physics-based interaction face many challenges including reasoning, planning, summarization, and question answering. This problem is exacerbated when a human user wishes to either guide or interact with the agent in natural language. Although the use of Language Models (LMs) is the default choice, as an AI tool, they struggle with tasks involving physics. The LM's capability for physical reasoning is learned from observational data, rather than being grounded in simulation. A common approach is to include simulation traces as context, but this suffers from poor scalability as simulation traces contain larger volumes of fine-grained numerical and semantic data. In this paper, we propose a natural language guided method to discover coarse-grained patterns (e.g., 'rigid-body collision', 'stable support', etc.) from detailed simulation logs. Specifically, we synthesize programs that operate on simulation logs and map them to a series of high level activated patterns. We show, through two physics benchmarks, that this annotated representation of the simulation log is more amenable to natural language reasoning about physical systems. We demonstrate how this method enables LMs to generate effective reward programs from goals specified in natural language, which may be used within the context of planning or supervised learning.

CCS Concepts: • **Computing methodologies** → **Artificial intelligence**; **Knowledge representation and reasoning**; *Natural language processing*; *Physical simulation*.

Additional Key Words and Phrases: Reasoning, Representation Learning

## 1 Introduction

Imagine a videogame designer using natural language to assess the feasibility of a task involving physical interaction. For example, consider a scene resembling Figure 1, where the designer wishes to know if it is possible to "Make the green ball bounce against the wall, bounce onto the table and land in the first bucket". This may require iterative simulations to find a matching trajectory. It is possible to detect a 'match' automatically, but the corresponding

Authors' Contact Information: Sean Memery, s.memery@ed.ac.uk, University of Edinburgh, Edinburgh, United Kingdom; Kartic Subr, k.subr@ed.ac.uk, University of Edinburgh, Edinburgh, United Kingdom.

program would take arduous crafting, where interaction tests would need to be performed on the input video (or state information) and tuned to match the conditions in the goal statement.

An obvious tool to exploit in such problems, involving natural language inputs, is a Large Language Model (LLM). While current Artificial Intelligence systems (AI) excel at interpreting the goal and even generating plans, they struggle to reliably interpret low-level simulation traces and dynamics [Liu et al. 2022; Mecattaf et al. 2024; Memery et al. 2025; Xu et al. 2025a]. Video and multimodal LMs exhibit limited success on intuitive physics benchmarks [Bordes et al. 2025; Memery et al. 2024; Shivan Jassim 2023; Xiang et al. 2025b]. These foundation models remain error-prone and often opaque, offering limited interpretability and explainability beyond the model's own textual justification [Kambhampati et al. 2024; Memery et al. 2025]. Natural language interaction has become a common technique in the graphics community in recent years, with many works utilising language based AI for guiding content generation, LM reasoning, and improved user interaction [Chen et al. 2025; Gao et al. 2023; Goel et al. 2024; Ji et al. 2024; Ma and Agrawala 2025; Menapace et al. 2024; Sun et al. 2024; Zhang et al. 2023].

We propose a method for extracting high-level patterns from raw simulation traces. Given a set of text snippets (e.g., "a lever launches a ball") describing potential patterns (say by a domain expert such as the game designer in the above example), we automatically synthesize programs that detect pattern activations from detailed simulation traces. We demonstrate that LLMs are able to effectively interpret these high-level annotated traces to generate summaries, perform question-and-answering (Q&A) tasks and to automatically synthesize formal and interpretable reward functions from natural language goals.

Given a simulation environment, we learn a library of pattern-detecting programs specialized to scenes within the environment. We use the Phyre benchmark [Bakhtin et al. 2019], as our 2D rigid-body environment. Seeded by natural language descriptions for
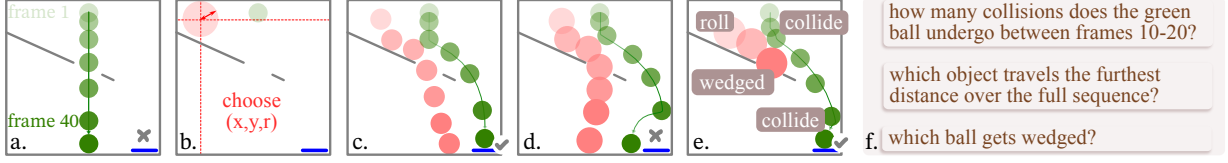
Fig. 2. (a–d) Example Phyre task. The objective for all Phyre tasks is to place the red ball so that the green and blue objects are in contact at the end of the simulation. (e) Highlighted patterns discovered by our system. (f) Example questions from our Q&A benchmark designed to probe physical reasoning.



Fig. 3. (a) Simulation traces $\tau_1, \tau_2$ are mapped to annotated simulation traces (ASTs) $A_1, A_2$ using detector code. Distance metrics $d_x$ and $d_p$ are defined between traces and ASTs. (b) We use FunSearch [Romera-Paredes et al. 2024], with a custom evaluation function, to augment the library with new detector code (c) Given a custom Domain Specific Language, a description of objects in the scene and the current library of pattern-detecting code, we synthesize a reward program in the DSL, which can be optimized to produce actions. Simulation traces produced by optimized actions are processed for reward evaluation.

potential patterns within the environment (roll, collide, bounce, etc.) we build a library of patterns, learning programs that can detect each pattern (by extending FunSearch [Romera-Paredes et al. 2024]). We use this library of patterns to translate low-level simulation trace data into *Annotated Simulation Traces* (AST) containing high level patterns identified by natural language descriptions. We show that, with ASTs as input, LLMs are effective at summarization, solving tasks, answering questions and synthesizing reward programs from natural language goals. Our ontributions are:

- we show the benefit of extracting high-level patterns from simulation traces, for applications involving natural language reasoning within environments with physical interaction;
- we invent a method that relies on minimal user guidance (a list of text snippets) to discover patterns from simulation data;
- we use evolutionary program synthesis to learn pattern detector programs to translate simulation traces into annotated traces;
- we demonstrate applications such as problem solving, Q&A, summarization and reward program synthesis.

## 2 Related work

### 2.1 Language models and physics environments

An ongoing thread of investigation in computational cognitive science is whether human intuitive physics is driven by approximate internal simulators or "intuitive physics engines" that support fast counterfactual prediction under uncertainty [Battaglia et al. 2013]. This assumes particular importance for interactive AI agents operating in physics environments. Verbal interaction and reasoning about physics naturally suggests the involvement of language models (LMs) and vision-language models (VLMs).

Recent work evaluates whether modern LMs and VLMs acquire comparable physical priors from large-scale pretraining, using both text-centric and video-centric benchmarks. Curated physics problem sets involving text and/or visual data as context is prominent [Chow et al. 2025; Dai et al. 2025; Qiu et al. 2025; Xiang et al. 2025a; Xu

et al. 2025b; Zhang et al. 2025]. However simulation-based benchmarks probe intuitive physical principles under temporal dynamics, some examples include GRASP [Shivan Jassim 2023] and Int-Phys2 [Bordes et al. 2025], which evaluate vision-language models intutive physics ability. Complementary evaluation suites target agentic performance in dynamic environments and games (e.g., BALROG [Paglieri et al. 2024]) and embodied interaction benchmarks (e.g., LLM-AAI [Mecattaf et al. 2024]) as well as explicitly physics-focused agentic VLM evaluation (e.g., DeepPHY [Xu et al. 2025a]). Across these settings, results generally indicate that strong language-based reasoning does not translate into robust physical prediction and control, and rarely approaches human level.

A common response is to ground reasoning in external tools or simulators. Mind's Eye [Liu et al. 2022] conditions an LM on outcomes generated by a physics simulation to improve physics question answering. It has also been shown that simulation can be used in closed-loop to ground LM reasoning [Cherian et al. 2024; Memery et al. 2025]. Several works argue that LMs are most reliable when paired with external model-based verifiers rather than used as standalone planners [Kambhampati et al. 2024; Memery et al. 2025, 2024]. In parallel, there is evidence that improving visual representations and training signals may be necessary for physics understanding. This includes the V-JEPA series of models [Assran et al. 2025; Garrido et al. 2025], where predicting physics outcomes within a learned representation space has shown promise. Additionally, reinforcement learning of vision-LMs in synthetic worlds has been shown to improve 3D embodied behavior [Bredis et al. 2025].

We build on two recent insights: First, that LMs are more effective when reasoning about high-level events rather than low-level simulation state traces [Memery et al. 2025, 2024]; and second, that LMs are effective at generating executable code that models environment dynamics and structure [Dainese et al. 2024; Romera-Paredes et al. 2024; Tang et al. 2024]. The latter learns executable transition models from interaction data. Inspired by these works, we learn

code to detect high-level patterns from simulation traces, and use these patterns to support LM reasoning about physical systems.

## 2.2 Reward program synthesis

Classical approaches to inverse reinforcement learning aim to infer rewards explaining expert behavior [Ng and Russell 2000]. However, we are interested in methods to recover structure and interpretibility, such as learning high-level task specifications with temporal structure [Vazquez-Chanlatte et al. 2018] or symbolic reward representations such as reward machines [Toro Icarte et al. 2018, 2019]. We are inspired by approaches that formulate reward learning as a program synthesis problem, inducing interpretable reward programs by example [Zhou and Li 2022], demonstrations or preferences.

Reward programs can be optimized to produce controllable behavior [Davidson et al. 2025; Yu et al. 2023] and are closely related to works that model environment dynamics and policies as programs for model-based reasoning and planning with strong compositional structure [Ahmed et al. 2025; Curtis et al. 2025; Tang et al. 2024]. Probabilistic-programming perspectives treat cognition and reasoning as the synthesis of task-specific generative models. For example, open-world cognition may be modeled as iterative construction and refinement of probabilistic models Wong et al. [2025]. Across these, a common theme is to use program synthesis to make goals, models, and evaluators explicit, allowing systems to adapt to new tasks while retaining interpretability and the ability to incorporate external verification. Across these approaches, LMs have been shown to be effective at generating candidate proposals for reward programs.

## 2.3 Program synthesis via FunSearch

FunSearch [Romera-Paredes et al. 2024], is a method for genetic programming (GP) and hence a form of evolutionary algorithm (EA). It retains the general loop structure of EAs, which is to maintain a population of candidate solutions and iteratively apply variation and selection to evolve better solutions. It searches the functional space of executable programs by replacing traditional rule-based or stochastic changes with a large language model (LLM) to handle mutation and discovery. Hallucination is controlled by an execution-based evaluation function that scores candidates within domain-specific contexts. It employs an 'island model' where multiple sub-populations evolve in parallel with occasional migration of high-performing candidates between islands to maintain diversity. Related methods incorporate evolutionary search or explicit reflective feedback to improve sample efficiency and exploration, including Evolution of Heuristics (EoH) [Liu et al. 2024] and ReEvo [Ye et al. 2024]. Tree-search variants such as GIF-MCTS [Dainese et al. 2024] further combine LM proposals with structured exploration to generate reliable code for environment modeling and planning.

We adapt this general scheme (see the algorithm in appendix F) to synthesize pattern detectors. Rather than optimizing towards a single labeled output, we score candidate programs by whether their emitted event streams covary with meaningful differences in trajectory geometry, while discouraging redundancy with respect to the current library (details in Sec. 3.2). We achieve this by providing (i) a user-defined evaluation function evaluate(·) to score candidate outputs (and rejects invalid ones). See appendix F for LM usage with FunSearch.

We use the open-source vision-language model Qwen3-VL 8B (Thinking) [Bai et al. 2025] and vLLM [Kwon et al. 2023] as an inference backend. For the code evolution pipeline only, we use Qwen3-Coder 30B due to its stronger programming performance. Additionally, we compare different sized models of the Qwen3-VL range, and also Nvidia's Cosmos-2 Reasoning model [NVIDIA 2025], which is trained on a Qwen3-VL backbone.

## 3 Method

### 3.1 Definitions: Patterns, annotation and detectors

Let $\mathbf{x}_i \in \mathcal{X}$ be the state in the state space of the physics environment at the $i^{th}$ time-step of the simulation and $\tau = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N\} \in \Upsilon$ be a simulation trajectory of length $|\tau| = N$. Let $\tau_i$ denote the $i^{th}$ trajectory while $\tau^{(i)}$ represents the $i^{th}$ step of $\tau$ ($\mathbf{x}_i$ in this case). We define an alternative, abstract space for simulation trajectories to enable high-level reasoning involving common events.

A *pattern* $\mathbf{p} \in \mathcal{P}$ captures a specific evolution of states within $\tau$. For example, a subsequence of states representing elastic collision. These subsequences are not mutually exclusive across states and so $\tau$ cannot be strictly defined as a sequence of patterns. We use a $|\tau| \times |\mathcal{P}|$ sparse binary *annotation matrix*, $\mathcal{P}$ is the *pattern library* and $A_{ij} = 1$ iff the $j^{th}$ pattern $\mathbf{p}_j$ is active at time-step $i$.

For each pattern $\mathbf{p}_j$ in the library, we define a *pattern detector* as a program that acts as a function $f_j : \Upsilon \times \Theta \to [0, 1]^N$. In the above example, $f_j(\tau, \theta)$ outputs the $j^{th}$ column of the annotation matrix, $A_{\cdot j}$ where $\theta_j \in \Theta$ represents pattern-specific parameters. For example, $\theta$ could contain identification numbers of objects involved in the pattern. Also associated with each pattern is a label $L_j$ which is a short natural language description of the pattern (e.g., "elastic collision between objects X and Y"), and a longer description $D_j$, which is used for more detailed explanations. Each pattern detector $f_j$ has access to the matrix $A$ as it is constructed, allowing it to reference activations of other patterns when determining its own activations. This creates a dependency graph among patterns, enabling the construction of complex patterns from simpler ones. See appendix B for a breakdown of the pattern detector dependencies.

We define two distance metrics to compare trajectories in the state and annotation space respectively. We define *trace distance* $d_{\mathbf{x}}$, as a normalized translational distance of matched objects. That is, if $O$ is the intersection of the sets of objects present in traces $\tau_1$ and $\tau_2$ respectively, then $d_{\mathbf{x}}(\tau_1, \tau_2)$ measures the average Euclidean distance between each object in $O$ per frame, across $\tau_1$ and $\tau_2$, normalized by the length of the trace.

We define *pattern annotation distance* $d_{\mathbf{p}}$, as the cross entropy between normalized histograms, of pattern occurrences over time, averaged across all patterns in the library. That is, given annotation matrices $A_1$ and $A_2$ and a pattern library with $J$ patterns, we discretize each into $b$ bins and compute normalized histogram counts $\mathbf{h}_j(A_1, b)$ and $\mathbf{h}_j(A_2, b)$ for $j = 1, 2, \cdots, J$. Then, $d_{\mathbf{p}}(A_1, A_2)$ is the cross entropy between $\mathbf{h}_j(A_1, b)$ and $\mathbf{h}_j(A_2, b)$, averaged across $J$ patterns. This distance captures how similar the distributions of activations of patterns across annotations of the two traces.

## 3.2 Natural language guided pattern discovery

Given a physics environment and a set of trajectories $\{\tau_k\}_{k=1}^K$, we learn a library of pattern detectors $\mathcal{P} = \{f_j\}_{j=1}^J$ that discover patterns from traces of simulation states $\mathbf{x}_1, \ldots, \mathbf{x}_N$. We construct a fitness function that scores candidate pattern detectors based on (1) how well the patterns they detect correlate with differences in trajectory geometry, and (2) how much new information they provide relative to patterns already in the library. We include penalty terms to discourage long programs, degenerate patterns and long execution time (details in Appendix F).

Our primary consideration is that patterns should reflect similarities between simulation traces. That is for two traces $\tau_1$ and $\tau_2$, with corresponding annotation matrices $A_1$ and $A_2$, if the traces are similar then their pattern annotations should also be similar, and vice versa. In other words, a high correlation between the distances in those spaces, $d_\mathbf{x}(\tau_1, \tau_2)$ and $d_\mathbf{p}(A_1, A_2)$, is desirable. In addition, we seek to discover patterns that are informative with respect to the existing library. For a candidate pattern detector $f_{new}$ producing annotation matrix $A_{new}$, we want $d_\mathbf{p}(f_j, f_{new})$ to be high, encouraging novelty. This is achieved by incorporating it in the fitness function ($\eta$ in Alg. 2)

We start from a pool of candidate pattern labels (natural language) provided by a user and use an evolutionary programming approach to search for corresponding pattern detectors that maximize the above fitness criteria. The pattern discovery algorithm (Algorithm 1) takes as input a set of trajectories $\mathcal{T}$, a set of candidate pattern labels $\mathcal{L}$ and a skeleton (seed) program $g_0$ which contains empty logic with the structure required of a pattern detector. After initializing $\mathcal{P}$ and seting up some parameters, it invokes FunSearch to synthesize a candidate pattern detector $f_m^*$ for each $L_m \in \mathcal{L}$ (along with its fitness score $v_m^*$). If $v_m$ exceeds a predefined threshold $\delta$, the candidate pattern detector is added to the library $\mathcal{P}$.

We prompt the language models with a structured prompt denoted $P_r(G, \mathcal{P}, \Theta, D, \{(\hat{G}_k, \hat{r}_k)\})$ where $G$ is the natural language goal , $\mathcal{P}$ is the pattern library with associated parameter sets $\Theta$, $D$ is a description of the DSL syntax and semantics, and $\{(\hat{G}_k, \hat{r}_k)\}$ is a set of few-shot examples of natural language goals and their corresponding reward programs (details in Appendix E). Since the synthesized programs may contain syntax errors, invalid identifier usage or mismatched parameter keys, if parsing or execution fails we iteratively prompt for automatic repair by supplying the candidate DSL program and the interpreter error message. We abort as a failure after a fixed retry limit is reached. See Appendix B for the full list of learned patterns, and Appendix E for LM prompts. FunSearch in Algorithm 1 Evaluate (Algorithm 2) as the fitness function for candidate pattern detectors. Given a set of trajectories $\mathcal{T}$, a candidate pattern detector $f_{new}$ and the current library $\mathcal{P}$, it computes the trace distances $d_\mathbf{x}$ and pattern annotation distances $d_\mathbf{p}$ for all pairs of trajectories in $\mathcal{T}$ using the current library and the candidate pattern detector. It also computes distances between annotations by the new pattern and annotations by the existing library, where a higher mean distance indicates greater novelty. Finally, it computes the correlation $\rho$ between $D_\mathbf{x}$ and $D_\mathbf{p}$, novelty score $\eta$, length penalty $\lambda$ and time penalty $\psi$, and combines them to produce the final fitness score $v$. Parameters $\theta$ for each pattern

---

**ALGORITHM 1:** DiscoverPatternDetectors

**Input:** Set of trajectories $\mathcal{T} = \{\tau_k\}_{k=1}^K$,
candidate pattern labels $\mathcal{L} = \{L_m\}_{m=1}^M$.
skeleton algorithm $g_0$
**Output:** Pattern library $\mathcal{P} = \{f_j\}_{j=1}^J$
Initialize empty pattern library $\mathcal{P} \leftarrow \{\}$;
Initialize LLM for FunSearch LLM($\cdot$);
Initialize FunSearch parameters $I, s, T_r$ ;
**for** *each label* $L_m$ *in* $\mathcal{L}$ **do**
$\quad (f_m^*, v_m^*) \leftarrow$ FunSearch(Evaluate($\mathcal{T}, \cdot, \mathcal{P}$), $g_0$, LLM, $I, s, T_r$);
$\quad$ **if** $v_m^* > \delta$ **then**
$\quad\quad \mathcal{P} \leftarrow \mathcal{P} \cup \{f_m^*\}$;

---

detector are inferred during synthesis and output as metadata by the synthesized program.

## 3.3 Reward program synthesis

Given a natural language goal (e.g., "make the red ball collide with the green object"), we synthesize a compositional expression in a custom domain-specific language (DSL) and call it a *reward program*. The reward program operates on an *annotated simulation trace* (AST) to produce (1) a boolean success/failure signal and (2) a dense reward signal in $[0, 1]$ indicating partial credit towards goal completion. The program can then be used as a reward function for trajectory optimization or supervised machine learning. The reward program $r(\cdot)$ is structured as a single DSL expression composed of multiple boolean predicates. When executed on an AST containing tuples of the form $(L_j, i, \theta_j)$ (where $L_j$ is a pattern label active at time-step $i$ with parameters $\theta_j$), $r$ serves as a test for whether the natural language goal $G$ was achieved. We use three classes of predicates and one quantitative primitive. *Event predicates* check for the occurrence of specific patterns within a trace. *Logical predicates* facilitate classical boolean operators such as AND, OR and NOT. *Temporal predicates* test activation timings and relative ordering of patterns in the trace. *Spatial and frequency quantifiers* measure spatial proximity and frequencies.

*Partial-credit scoring.* In addition to boolean satisfaction, we compute a dense reward in $[0, 1]$. We interpret the synthesized program to be composed of a top-level AND operator with multiple operands and return a reward that is the average of the number of subclauses that evaluate to true. For quantitative primitives, we assign graded scores based on the distance-to-satisfaction. For NEARBY_AT, we convert the object-to-target distance into a score using an inverse-log transform and clamp to $[0, 1]$, so improvements near the target are weighted strongly. For COUNT and comparisons, we compute a deviation from the target count and map it to $[0, 1]$ again with an inverse-log shaping and clamping. For example, the goal "Launch the green ball into the second bucket", for the Phyre scene shown in Figure 1(d), may be synthesized into the following reward program:

```
AND(
  # Check for lever launch event
  EVENT("lever launch", {"object": "green ball"}),

  # Check for collision with object after launch event
  AFTER("lever_launch", "collision",
        first_params={"object": "green ball"},
```

```
            second_params={"object_a": "green ball",
                           "object_b": "object X"}),

  # Green ball follows a particular path of grid cells
  EVENT("grid cell path", {"object": "green ball",
        "path": [(5,10), (10,10), (12,5), (14,2)]}),

  # Check if green ball is inside target bucket at end
  NEARBY_AT("green ball", x=140, y=10, t=1.0),
)
```

Here, we check for the existance of an "lever_launch" event involving the green ball. The "AFTER" call enforces that the "collision" event involving the green ball and the wall must occur after the launch event. The "grid cell path" event specifies the exact path the green ball follows through the grid. Finally, the NEARBY_AT predicate checks if the green ball is near the specified coordinates (where the scene is divided into a $25 \times 25$ grid) at the end of the trajectory.

---

**ALGORITHM 2:** Evaluate

**Input:** Set of trajectories $\mathcal{T} = \{\tau_k\}_{k=1}^{K}$,
         Candidate pattern detector $f_{\text{new}}$,
         Current library $\mathcal{P}$

**Output:** Fitness score $\nu$

Initialize empty lists $D_{\mathbf{x}} \leftarrow [\,], D_{\mathbf{p}} \leftarrow [\,], D_{\text{novel}} \leftarrow [\,]$;

**for** *each pair of trajectories* $(\tau_l, \tau_m)$ *in* $\mathcal{T}$ **do**
     Compute annotation matrices $A_l, A_m$ using current library $\mathcal{P}$;
     Compute annotation vectors $\mathbf{a}_l, \mathbf{a}_m$ using $f_{\text{new}}$;
     Append $d_{\mathbf{x}}(\tau_l, \tau_m)$ to $D_{\mathbf{x}}$;
     Append $d_{\mathbf{p}}(\mathbf{a}_l, \mathbf{a}_m)$ to $D_{\mathbf{p}}$;
     Append $d_{\mathbf{p}}(\mathbf{a}_l, A_l)$ and $d_{\mathbf{p}}(\mathbf{a}_m, A_m)$ to $D_{\text{novel}}$;

$\rho \leftarrow \text{corr}(D_{\mathbf{x}}, D_{\mathbf{p}})$;

$\eta \leftarrow \text{mean}(D_{\text{novel}})$;

$\lambda \leftarrow \text{ComputeLengthPenalty}(f_{\text{new}})$;

$\psi \leftarrow \text{ComputeTimePenalty}(f_{\text{new}})$;

Compute fitness score: $\nu \leftarrow \rho + \eta - \lambda - \psi$;

---

## 4 Evaluation of discovered patterns

### 4.1 Datasets and patterns

We evaluate $\mathcal{P}$ on two benchmarks involving physical reasoning in 2D rigid-body simulations: an adapted Phyre task suite, to solve puzzles [Bakhtin et al. 2019], and a question answering (Q&A) benchmark constructed from Phyre rollouts. In both, LMs are provided with basic initial state information: object IDs, colour, shape, and position. The LMs are not given $\tau$ explicitly (i.e. the low-level simulation trace data), and are shown only video data or annotations.

*Phyre task suite.* The dataset consists of 2D rigid-body puzzles where an agent performs a single action, placing a red circle with parameters $(x, y, r)$, where $x, y \in [0, 256]$ and $r \in [4, 32]$, to induce physical interactions that terminate with the green and blue objects in contact by the end of the rollout. We restrict our experiments to the 25 one_ball templates from the original benchmark, each with 100 scene variants, yielding 2500 tasks. We execute tasks using a PyGame-based reimplementation to record per-timestep state and event traces. To make large scale experimentation with LMs feasible,

we quantize the action space using $4 \times 4 \times 3$ bins, yielding 48 possible choices for actions, per task. We preprocess each task, solving it via random choices for actions, to ensure that there exists a solution action. For LM interaction, we provide a natural-language task description and the enumerated action choices. We additionally render the initial scene (and subsequent rollout videos when operating in multi-attempt settings) for video capable LMs [Bai et al. 2025].

*Q&A benchmark.* We construct a new benchmark from Phyre rollouts. We sample 100 scene instances and, for each instance, select an action by sampling uniformly between: a stored successful solution, or a near-miss action obtained by perturbing a solution in action space. For each resulting rollout we generate 10 questions by instantiating templates that target counts and temporal statistics, categorical identification, and prediction under partial observation. Example numerical questions include: "How many objects did the red ball touch between times $t = 0.5s$ and $t = 1.5s$?" Example predictive questions include: "Will the green object touch the blue object after $t = 2.0s$?" where the model is only given the rollout video and, when applicable, annotations, up to $t = 2.0s$. Answers are computed deterministically from the recorded trace. Full question templates are provided in the Appendix C.

*Guided and self-discovered patterns.* In practice, once detectors have been synthesized for all labels provided by the domain expert, we begin to generate new labels automatically. We provide the LM with (i) the existing set of patterns $\mathcal{P}$ along with their natural-language descriptions and (ii) reasoning traces collected from responses on the Q&A benchmark. This method allows the LM to find common patterns across different tasks and scenarios. This leads to a phase of discovery where new labels are discovered along with their corresponding pattern detectors. We then evaluate newly discovered patterns with the Q&A benchmark, and only retain those that improve performance. This evaluation then fuels further rounds of pattern discovery. We refer to these two subsets of patterns as *guided* and *automatic* respectively. The aim is to encourage the discovery of automatic patterns that build on existing guided patterns enabling higher levels of abstraction. The final output is the learned pattern library containing all discovered pattern detectors.

### 4.2 Size of the pattern library

We assessed the impact of the size of the pattern library on performance on our benchmarks. We sample sets of patterns of increasing size and with each set we evaluate performance on the Q&A and Phyre benchmarks using annotations generated by the sampled detectors. We averaged performance over 25 trials for each datapoint. Importantly, for the Phyre benchmark, we max the number of attempts to 5 per task, to stay in our computational budget.

Figure 4 shows how performance improves with library size. While performance on both benchmarks increase steadily until about 10 patterns, it pleateaus at around 12. Beyond 16, we include all guided patterns and begin to add automatic patterns. We believe that the flattening of performance could be related to the complexity of the environment and the dataset (quantization of action space, the complexity of questions, etc.) as well as to the nature of patterns.
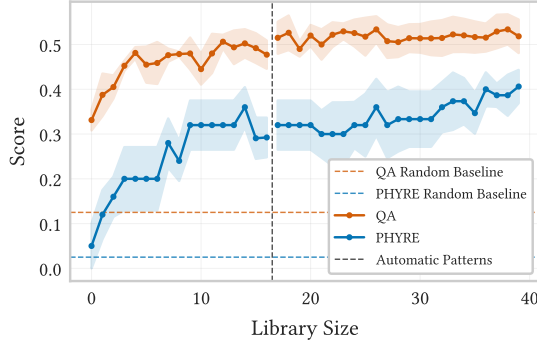
Fig. 4. Effect of library size on performance within the Q&A and Phyre benchmarks, using the Qwen3-VL 8B (Thinking) language model. Beyond $|\mathcal{P}| = 16$, automatic patterns were included. Not that the number of attempts at the Phyre task is capped at 5 per scene.
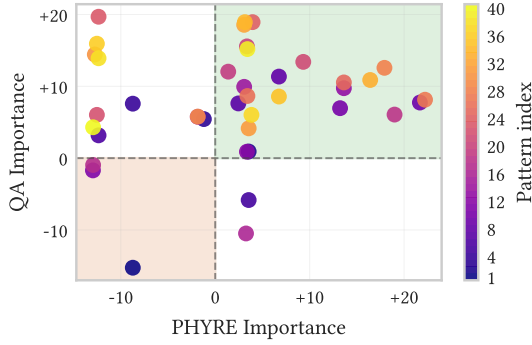


Fig. 5. Importance values for each pattern the learned library, computed via leave-one-out ablation on the Q&A and Phyre benchmarks. Importance is measured as the percentage drop in performance when ablating each pattern (higher is more important).

For example, one particular guided label is "An object that was not directly part of a collision begins to move shortly after the collision occurs, indicating it was supported or constrained by one of the colliding objects". This pattern captures the important physical concept of indirect contact via support, which is critical for reasoning about object interactions in Phyre tasks. In contrast, an automatic label is "The centers of the green and blue objects achieve near-equal x or y coordinates (within a small tolerance), indicating a horizontal or vertical axis alignment for the first time in the run". While this pattern captures a geometric relationship, it is less directly relevant to the physical interactions that are central to the tasks. Another automatic label, "Green and blue make contact and later separate before the final frame", also does not capture an important physical relationship, but it does directly relate to the task goal of bringing the green and blue objects into contact. Thus, human expertise could guide the selection of patterns relevant to the domain or tasks.



Fig. 6. Example annotations produced by the learned library on Phyre rollouts. Each row shows a different scene, with rendered frames on the left and time-stamped pattern activations on the right.

### 4.3 Illustrative examples

Figure 6 shows example annotations produced by pattern detectors from the learned library on Phyre rollouts. Each row contains a different scene, with simulation frames on the left and the corresponding pattern occurrences (with frame timestamps) on the right. The annotations highlight key interactions between objects and when they occur. For example, in the third row the library captures a lever interaction that launches the red ball, which is difficult to identify reliably from raw state traces alone.

Discrete pattern annotations also help with grounding. Pattern names make it easy to refer to the same interaction when it happens multiple times, and pattern parameters provide consistent references to the objects involved. This is especially useful in multi-object

Green ball rolls off the left side of the platform, landing in the corner.

```
AND(
    EVENT("collide", "green-ball", "red-ball"),
    AFTER("lose support", "collide",
            "green-ball"),
    EVENT("collide", "green-ball", "left-wall")
)
```

Red and green balls roll onto the blue bar, getting wedged into place.

```
AND(
    EVENT("roll over", "green-ball",
            "black-bar"),
    EVENT("roll over", "red-ball", "black-bar"),
    EVENT("support", "blue-bar",
            "green-ball"),
    AFTER("support", "wedged", "red-ball",
            "black-bar", "green-ball"),
)
```

The red ball causes the green ball to pass between the black objects without colliding with more than two.

```
AND(
    EVENT("collide", "red-ball", "green-ball"),
    AFTER("collide",
            NEARBY_AT("green-ball",
            x=180,y=100)
    ),
    LT(COUNT("collide", "green-ball",
            "black-ball"), 3)
)
```

Make the green and blue balls are touching at rest but make sure the red ball does not touch the green ball.

```
AND(
    NOT( EVENT("collide", "green-ball", "red-ball") ),
    EVENT("collide", "blue-ball", "red-ball"),
    AFTER("collide",
            EVENT("cell transition pattern",
                "blue-ball", "west")
            EVENT("CollisionStart", "blue-ball", "green-ball" )
    ),
    NOT( EVENT( "CollisionEnd", "blue-ball",
            "green-ball" ) )
)
```
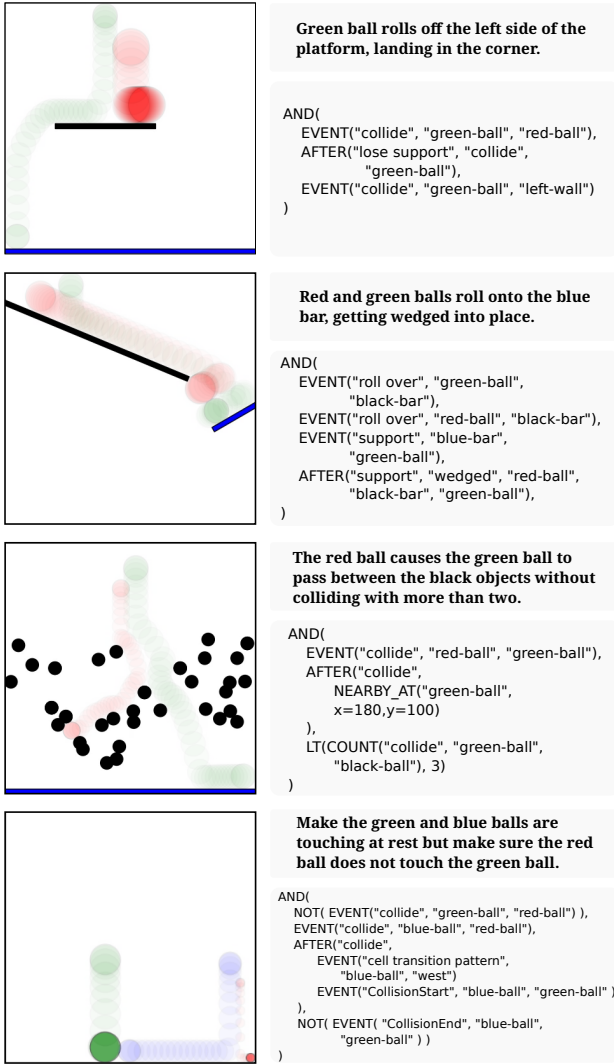
Fig. 7. Example reward programs generated by our system. On the left, the visualised rollout shows the optimised action for a Phyre task, with the natural language description of the task and the generated reward program shown on the right.

scenes, where LMs conditioned only on video often struggle to keep object references consistent over time.

Figure 7 shows example reward programs generated by our system for the Phyre benchmark. Each row corresponds to a different scene and task, with the rollout visualisation on the left and the natural language task description and generated reward program on the right. Our system generated the reward program in a simple DSL from the natural language task description, using the learned pattern library to define reward conditions. The generated reward program is then used to guide action selection by optimizing the reward program through simulated annealing.
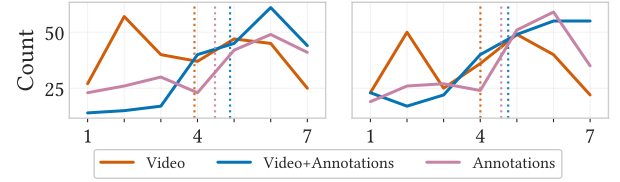


Fig. 8. Human evaluation results for LM-generated simulation trace summaries under input conditions: (i) Video only, (ii) Video and annotations of learned patterns, (iii) Annotations only. On the left are ratings of summaries from the guided library only, and on the right are ratings from the full library (guided + automatic). Dashed vertical lines indicate mean scores.

In each case, the generated reward program captures the key conditions for task success. For example, in the second row, the user provides a method of solving the Phyre task (with the standard success condition that the green and blue objects are touching at the end of the rollout) by having the red and green balls be wedged between the black and blue objects. Our learned library contains the relevant patterns to capture this method, in this case the event of the red ball being wedged between the black and green objects, while the green ball is touching the blue bar. The generated reward program then captures the necessary conditions for success, which the optimization process can use to find a successful action. In the third row, the task description adds an additional constraint that the green ball cannot collide with more than two of the black objects. This constraint must be understood by the LM and specified in the reward program, using the pattern library to define the relevant conditions and DSL constructs to express the logic.

## 4.4 Quantitative assessment of discovered patterns

To quantify the contributions of each pattern in $\mathcal{P}$, we perform a leave-one-out ablation over the learned libraries (natural-language guided and guidance-free) towards the downstream tasks of action selection (Phyre) and question answering (Q&A). First, we evaluate the performance of full library on each benchmark to obtain baselines. We then re-evaluate the performance $J$ times, each time removing a different $f_j$ and any other patterns that depend on it. We report the difference of the performance with each hold-out from the baseline, averaging over 100 runs with different random seeds. We define the *importance* of each pattern $f_j$ as the percentage drop in performance when it is held out. Figure 5 shows the importance values for each pattern on the Q&A and Phyre benchmarks. Certain patterns have a significant impact on performance. Automatic patterns (indices greater than 16) tend to exhibit less impact, despite many contributing positively in the Q&A benchmark.

## 5 Applications of discovered patterns

## 5.1 Pattern annotation helps summarization

We conducted a human study on the Prolific [2026] platform. Participants evaluated summaries generated by LMs that were provided three different inputs: video only, video plus pattern annotations and pattern annotations only. We used both the guided pattern library and the guided + automatic library for this experiment. We
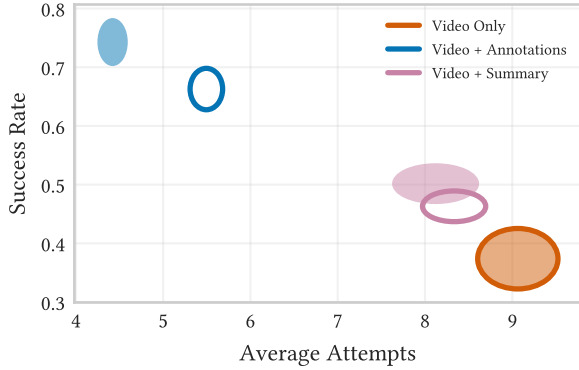
Fig. 9. Phyre task performance against the average number of attempts before success, using the Qwen3-VL 8B (Thinking) language model, for three variants of feedback: (i) Video only, (ii) Video and annotations of learned patterns, (iii) Video and an LM summary generated from the annotations. Results from the guided patterns only are shown as hollow ellipses, and results from the full library (guided + automatic) are shown as filled ellipses. Ellipses indicate standard errors over number of attempts (X-axis) and success rate (Y-axis). Not that the number of attempts at the Phyre task is capped at 15 per scene.
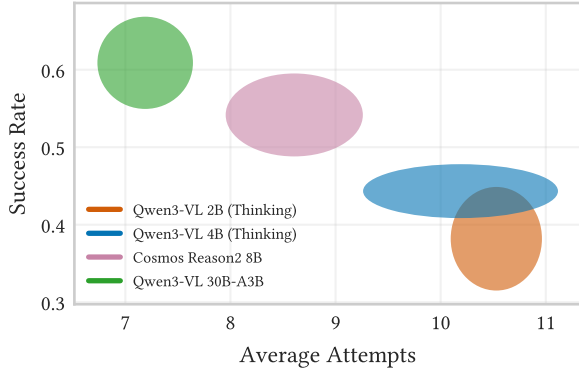


Fig. 10. Phyre task performance, against the average number of attempts before success, for video and annotations feedback only, across multiple LMs. Ellipses indicate standard errors over number of attempts (X-axis) and success rate (Y-axis).

used the Qwen3-VL 8B (Thinking) model to generate summaries. All 150 participants view pairs of rollout videos and single summaries. Participants rate each summary on a 1–7 likert scale based on how accurately it describes the video. The input used to generated the summary is sampled uniformly from the three conditions above. Full details and example summaries are in Appendix D.

Figure 8 plots histograms of the scores received by each summarization input. Using video and annotations yields the largest expected scores and received fewest low scores. Although the overall difference in mean scores (dashed vertical lines) is relatively small, indicating that human raters found all summaries to be reasonably accurate, the fact that reasonable summaries can be generated using annotations only, without video, indicates that the learned patterns capture salient physical events in the simulation.
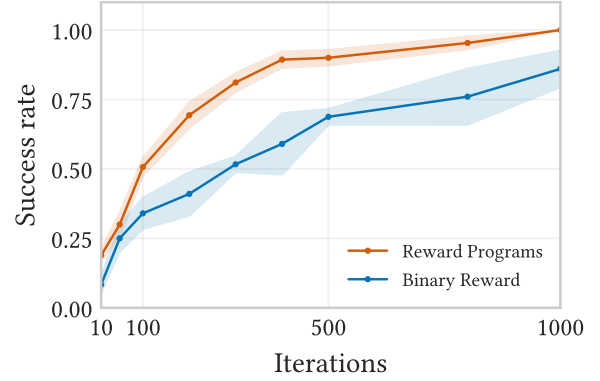


Fig. 11. Success rates for a single scene (see figure 1(d)) when optimising synthesized DSL rewards versus binary rewards with simulated annealing, across varying optimization budgets.

In addition to descriptive text, we measure downstream performance on the Phyre benchmark over the full set of 25 scenes, where an LM receives feedback after each action. We compared three variants: video only, video plus pattern annotation, and pattern annotation only. We measured the average task success rate for each variant over all scenes with 10 seeds each. Additionally, we recorded the average number of attempts made by the LM before it accomplished the task, where the max number of attempts allowed was 12. Figure 9 shows the performance of the two annotation based methods, compared to the video only baseline. The width and height of each ellipse depict standard errors over number of attempts and success rate respectively. Using video and annotations as feedback on actions yields the best performance. Using summaries showed only a slight improvement over video-only, with a decrease in performance compared to video plus annotations.

## 5.2 Achieving natural language goals via optimization

We optimized synthesized reward programs on scenes and measuring success rates. We use a single scene (Figure 1) with three goal prompts: "Get the green ball in the first bucket", "Get the green ball in the second bucket", and "Get the green ball in the third bucket". For each prompt, we synthesize a reward program and optimize the action parameters using simulated annealing. We tested with increasing annealing samples $N_o$ of 10, 50, 100, . . . , 1000. We run simulated annealing to select the highest-scoring candidate action under the reward program and execute it in the simulator. We count the trial as a success if the final distance between the green ball and the target bucket specified by the prompt is less than 10 pixels. We repeat this experiment over multiple random seeds and report the average success across $N_o$.

We compare the our synthesized reward programs, against sparse binary reward functions that return a reward of 1 if the goal is achieved (green ball in target bucket) and 0 otherwise, and we run the same simulated annealing procedure for both reward types. Figure 11 shows an improvement in optimization success rates when using synthesized reward programs versus binary rewards, across all optimization budgets. The curve for synthesized rewards rises
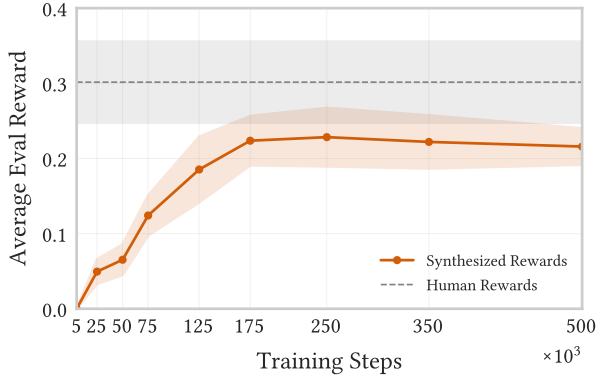
Fig. 12. Training steps (x-axis) versus evaluation performance (y-axis), measured as mean DSL reward-program score from suggested actions on 25 validation scenes, averaged across 10 goals. Five rewards created by hand are also shown for comparison, trained once at 250$k$ steps.

steeply, achieving over 75% success with only 250 optimization samples, while the binary reward curve rises slowly, achieving only 50% success with the equivalent optimization samples. This trend continues, showing the sample efficiency benefits of using synthesized reward programs that provide dense feedback during optimization.

## 5.3 Training value networks on reward programs

We used our reward programs to train value models to predict reward satisfaction, and used this to solve Phyre tasks. We constructed a trace dataset of 10000 rollouts sampled from the Phyre benchmark and annotate each trace using $\mathcal{P}$. Each scene is rendered as a 256×256 7-channel grid, where channels encode per-pixel roles: (1) dynamic goal object, (2) static goal subject, (3) dynamic goal subject, (4) static confounding body, (5) dynamic confounding body, (6) body placed by the agent, and (7) background. We define 10 human-written NL goals (see Appendix G) and synthesize one DSL reward program per goal. For each goal, we train an action-conditioned value model to predict reward satisfaction from the semantic grid and the proposed action, using a ResNet18 backbone with FiLM fusion of the action vector and balanced positive/negative training examples. We used the Adam optimizer, a learning rate of $10^{-3}$, and batch size of 64, matching the Phyre paper specification [Bakhtin et al. 2019]. During training, we periodically evaluate on 25 held-out validation scenes by uniformly sampling 1024 valid actions per scene, selecting the highest-scoring action under the value model, executing it in the simulator, and measuring the achieved reward score. We report evaluation performance as our mean reward score averaged across validation scenes and goals. We constructed an additional five reward functions by hand (see Appendix G for details). We repeated the training procedure above using these hand-written reward functions and use the maximum performance achieved after 250$k$ training steps as a baseline. Figure 12 shows that our reward programs lead to an increase in evaluation performance on unseen scenes. They are able to train value models, but fall short of human performance ( 0.24 vs 0.31 after 250 steps).

## 6 Discussion

*Importance of discovered patterns.* The negative importance of some patterns in Figure 5 raises questions about whether these should be ignored. For example, Pattern 2 is labeled "grid cell transition" and activates when an object transitions between cells on a 25 × 25 grid overlaid on the image. This is expected to trigger at high frequency during simulation, introducing clutter that could potentially confuse downstream reasoning. The need for mapping grid-cell identities back to locations in the video may exacerbate reasoning for the LM. However, Pattern 18 builds on Pattern 2 by summarizing sequences of cell transitions into coarse directions (e.g., "north-east", "south-west"). Interestingly, this higher-level abstraction has a substantially larger impact, and is the most beneficial pattern for Phyre performance. This is intriguing but we leave investigations into hierarchies of patterns for future work.

*Patterns as helpful abstractions.* By translating simulation data into annotations of high-level events, patterns enable LMs to focus on salient physical interactions, improving their ability to answer questions and select actions (Figure 4). However, not all patterns are equally beneficial, and some introduce noise if they trigger too frequently or lack semantic clarity (Figure 5). Yet such patterns could serve as stepping stones to the discovery of important patterns when combined with others (case in point, Pattern 18). A possible reason for the effectiveness of patterns to improve Phyre task performance more than Q&A accuracy is that patterns are designed to capture physical interactions, a vital aspect of Phyre tasks, whereas Q&A enquires about details that may not align with the patterns. Additionally, Phyre tasks are more difficult than Q&A questions, so the benefits of patterns may be more pronounced. The pattern library serves as a structured interface that bridges the gap between raw simulation data and LM reasoning capabilities and the resulting annotations complement video data to improve downstream tasks such as summarization and reward synthesis.

*Annotations enable effective reward synthesis.* By beating the sparse rewards in optimization (figure 11), we show that synthesized reward programs both (i) capture the intended goals specified in natural language, and (ii) provide dense feedback that supports sample-efficient optimization. Synthesized rewards can be easily adapted in natural language through LM refinement, making them a flexible tool for specifying goals. In the more general case, synthesized reward programs reach near-human average rewards (figure 12), demonstrating that they provide an effective training signal for value models. While there remains a gap to human performance, this is a promising step toward flexible goal specification through natural language, enabled by the learned pattern library.

*Limitations and future work.* Although our pattern discovery and annotations form an abstract, compressed and interpretable representation for generalized and grounded verbal reasoning about physics environments, we highlight a few of its limitations. First, we rely on some human expert labels (in natural language) for seeding pattern discovery. While our method is expected to generalize across environments, our demonstrations are limited to 2D rigid body interactions within Phyre. Empirical validation beyond this world would be necessary to assess its capability for generalization.

Our experiments focus on relatively simple tasks and therefore small libraries, up to 16 seed patterns, are sufficient. It remains to be seen whether larger libraries would result in longer, denser annotations which could pose challenges due to long contexts [Liu et al. 2023].

## 7 Conclusion

Our results demonstrate that learning an event-based pattern library from simulation traces provides a practical interface between physics environments and language models. Across question answering, summarization, and the Phyre benchmark, these annotations improve LM performance on interpreting and reasoning about physical rollouts. Finally, by grounding reward program synthesis in the learned pattern library, we enable executable goal specifications from natural language that support both optimization and supervised machine learning.

## References

Zergham Ahmed, Joshua B. Tenenbaum, Christopher J. Bates, and Samuel J. Gershman. 2025. Synthesizing world models for bilevel planning. arXiv:2503.20124 [cs.AI] https://arxiv.org/abs/2503.20124

Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Mojtaba, Komeili, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, Sergio Arnaud, Abha Gejji, Ada Martin, Francois Robert Hogan, Daniel Dugas, Piotr Bojanowski, Vasil Khalidov, Patrick Labatut, Francisco Massa, Marc Szafraniec, Kapil Krishnakumar, Yong Li, Xiaodong Ma, Sarath Chandar, Franziska Meier, Yann LeCun, Michael Rabbat, and Nicolas Ballas. 2025. V-JEPA 2: Self-Supervised Video Models Enable Understanding, Prediction and Planning. arXiv:2506.09985 [cs.AI] https://arxiv.org/abs/2506.09985

Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng Ren, Xingzhang Ren, Sibo Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang, Xi Zhang, Bo Zheng, Humen Zhou, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and Ke Zhu. 2025. Qwen3-VL Technical Report. arXiv:2511.21631 [cs.CV] https://arxiv.org/abs/2511.21631

Anton Bakhtin, Laurens van der Maaten, Justin Johnson, Laura Gustafson, and Ross Girshick. 2019. PHYRE: A New Benchmark for Physical Reasoning. arXiv:1908.05656 [cs.LG] https://arxiv.org/abs/1908.05656

Peter W. Battaglia, Jessica B. Hamrick, and Joshua B. Tenenbaum. 2013. Simulation as an Engine of Physical Scene Understanding. Proceedings of the National Academy of Sciences 110, 45 (Nov. 2013), 18327–18332. doi:10.1073/pnas.1306572110

Florian Bordes, Quentin Garrido, Justine T Kao, Adina Williams, Michael Rabbat, and Emmanuel Dupoux. 2025. IntPhys 2: Benchmarking Intuitive Physics Understanding In Complex Synthetic Environments. (2025). doi:10.48550/ARXIV.2506.09849

George Bredis, Stanislav Dereka, Viacheslav Sinii, Ruslan Rakhimov, and Daniil Gavrilov. 2025. Enhancing Vision-Language Model Training with Reinforcement Learning in Synthetic Worlds for Real-World Success. arXiv:2508.04280 [cs] doi:10.48550/arXiv.2508.04280

Bohong Chen, Yumeng Li, Youyi Zheng, Yao-Xiang Ding, and Kun Zhou. 2025. Motion-example-controlled Co-speech Gesture Generation Leveraging Large Language Models. In Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers (SIGGRAPH Conference Papers '25). Association for Computing Machinery, New York, NY, USA, Article 55, 12 pages. doi:10.1145/3721238.3730611

Anoop Cherian, Radu Corcodel, Siddarth Jain, and Diego Romeres. 2024. LLMPhy: Complex Physical Reasoning Using Large Language Models and World Models. (2024). doi:10.48550/ARXIV.2411.08027

Wei Chow, Jiageng Mao, Boyi Li, Daniel Seita, Vitor Campanholo Guizilini, and Yue Wang. 2025. PhysBench: Benchmarking and Enhancing Vision-Language Models for Physical World Understanding. ArXiv abs/2501.16411 (2025). https://api.semanticscholar.org/CorpusID:275932005

Aidan Curtis, Hao Tang, Thiago Veloso, Kevin Ellis, Joshua Tenenbaum, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. 2025. LLM-Guided Probabilistic Program Induction for POMDP Model Estimation. arXiv:2505.02216 [cs.AI] https://arxiv.org/abs/2505.02216

Song Dai, Yibo Yan, Jiamin Su, Dongfang Zihao, Yubo Gao, Yonghua Hei, Jungang Li, Junyan Zhang, Sicheng Tao, Zhuoran Gao, and Xuming Hu. 2025. PhysicsArena: The First Multimodal Physics Reasoning Benchmark Exploring Variable, Process, and Solution Dimensions. ArXiv abs/2505.15472 (2025). https://api.semanticscholar.org/CorpusID:278783124

Nicola Dainese, Matteo Merler, Minttu Alakuijala, and Pekka Marttinen. 2024. Generating Code World Models with Large Language Models Guided by Monte Carlo Tree Search. arXiv:2405.15383 [cs]

Guy Davidson, Graham Todd, Julian Togelius, Todd M. Gureckis, and Brenden M. Lake. 2025. Goals as reward-producing programs. Nature Machine Intelligence 7, 2 (Feb. 2025), 205–220. doi:10.1038/s42256-025-00981-4

William Gao, Noam Aigerman, Thibault Groueix, Vova Kim, and Rana Hanocka. 2023. TextDeformer: Geometry Manipulation using Text Guidance. In ACM SIGGRAPH 2023 Conference Proceedings (Los Angeles, CA, USA) (SIGGRAPH '23). Association for Computing Machinery, New York, NY, USA, Article 82, 11 pages. doi:10.1145/3588432.3591552

Quentin Garrido, Nicolas Ballas, Mahmoud Assran, Adrien Bardes, Laurent Najman, Michael Rabbat, Emmanuel Dupoux, and Yann LeCun. 2025. Intuitive Physics Understanding Emerges from Self-Supervised Pretraining on Natural Videos. (2025). doi:10.48550/ARXIV.2502.11831

Purvi Goel, Kuan-Chieh Wang, C. Karen Liu, and Kayvon Fatahalian. 2024. Iterative Motion Editing with Natural Language. In ACM SIGGRAPH 2024 Conference Papers (Denver, CO, USA) (SIGGRAPH '24). Association for Computing Machinery, New York, NY, USA, Article 71, 9 pages. doi:10.1145/3641519.3657447

Xuebo Ji, Zherong Pan, Xifeng Gao, and Jia Pan. 2024. Text-Guided Synthesis of Crowd Animation. In ACM SIGGRAPH 2024 Conference Papers (Denver, CO, USA) (SIGGRAPH '24). Association for Computing Machinery, New York, NY, USA, Article 105, 11 pages. doi:10.1145/3641519.3657516

Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri, Lucas Saldyt, and Anil Murthy. 2024. LLMs Can't Plan, But Can Help Planning in LLM-Modulo Frameworks. (2024). doi:10.48550/ARXIV.2402.01817

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. arXiv:2309.06180 [cs.LG] https://arxiv.org/abs/2309.06180

Fei Liu, Xialiang Tong, Mingxuan Yuan, Xi Lin, Fu Luo, Zhenkun Wang, Zhichao Lu, and Qingfu Zhang. 2024. Evolution of Heuristics: Towards Efficient Automatic Algorithm Design Using Large Language Model. arXiv:2401.02051 [cs]

Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the Middle: How Language Models Use Long Contexts. arXiv:2307.03172 [cs.CL] https://arxiv.org/abs/2307.03172

Ruibo Liu, Jason Wei, Shixiang Shane Gu, Te-Yen Wu, Soroush Vosoughi, Claire Cui, Denny Zhou, and Andrew M. Dai. 2022. Mind's Eye: Grounded Language Model Reasoning through Simulation. (2022). doi:10.48550/ARXIV.2210.05359

Jiaju Ma and Maneesh Agrawala. 2025. MoVer: Motion Verification for Motion Graphics Animations. ACM Trans. Graph. 44, 4, Article 33 (July 2025), 17 pages. doi:10.1145/3731209

Matteo G. Mecattaf, Ben Slater, Marko Tešić, Jonathan Prunty, Konstantinos Voudouris, and Lucy G. Cheke. 2024. A Little Less Conversation, a Little More Action, Please: Investigating the Physical Common-Sense of LLMs in a 3D Embodied Environment. (2024). doi:10.48550/ARXIV.2410.23242

Sean Memery, Kevin Denamganaï, Jiaxin Zhang, Zehai Tu, Yiwen Guo, and Kartic Subr. 2025. CueTip: An Interactive and Explainable Physics-aware Pool Assistant. In Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers (SIGGRAPH Conference Papers '25). Association for Computing Machinery, New York, NY, USA, Article 86, 11 pages. doi:10.1145/3721238.3730742

Sean Memery, Mirella Lapata, and Kartic Subr. 2024. SimLM: Can Language Models Infer Parameters of Physical Systems? arXiv:2312.14215 [cs.CL] https://arxiv.org/abs/2312.14215

Willi Menapace, Aliaksandr Siarohin, Stéphane Lathuilière, Panos Achlioptas, Vladislav Golyanik, Sergey Tulyakov, and Elisa Ricci. 2024. Promptable Game Models: Text-guided Game Simulation via Masked Diffusion Models. ACM Trans. Graph. 43, 2, Article 17 (Jan. 2024), 16 pages. doi:10.1145/3635705

Andrew Y. Ng and Stuart Russell. 2000. Algorithms for Inverse Reinforcement Learning. Proceedings of the Seventeenth International Conference on Machine Learning (2000).

NVIDIA. 2025. cosmos-reason2. https://github.com/nvidia-cosmos/cosmos-reason2. GitHub repository. Accessed: 2026-01-13.

Davide Paglieri, Bartłomiej Cupiał, Samuel Coward, Ulyana Piterbarg, Maciej Wolczyk, Akbir Khan, Eduardo Pignatelli, Łukasz Kuciński, Lerrel Pinto, Rob Fergus, Jakob Nicolaus Foerster, Jack Parker-Holder, and Tim Rocktäschel. 2024. BALROG: Benchmarking Agentic LLM and VLM Reasoning On Games. arXiv:2411.13543

Prolific. 2026. Prolific. https://www.prolific.com/.

Shi Qiu, Shaoyang Guo, Zhuo-Yang Song, Yunbo Sun, Zeyu Cai, Jiashen Wei, Tianyu Luo, Yixuan Yin, Haoxu Zhang, Yi Hu, Chenyang Wang, Chencheng Tang, Haoling Chang, Qi Liu, Ziheng Zhou, Tianyu Zhang, Jingtian Zhang, Zhangyi Liu, Minghao

Li, Yuku Zhang, Boxuan Jing, Xianqi Yin, Yutong Ren, Zizhuo Fu, Weike Wang, Xu Tian, Anqi Lv, Laifu Man, Jianxiang Li, Feiyu Tao, Qihua Sun, Zhou Liang, Yu-Song Mu, Zhong-wei Li, Jing-Jun Zhang, Shutao Zhang, Xiaotian Li, Xingqi Xia, Jiawei Lin, Zheyu Shen, Jiahang Chen, Qi Xiong, Binran Wang, Fengyuan Wang, Ziyang Ni, Bohan Zhang, Fan Cui, Changkun Shao, Qing-Hong Cao, Mingjian Luo, Muhan Zhang, and Hua Xing Zhu. 2025. PHYBench: Holistic Evaluation of Physical Perception and Reasoning in Large Language Models.

Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M. Pawan Kumar, Emilien Dupont, Francisco J. R. Ruiz, Jordan S. Ellenberg, Pengming Wang, Omar Fawzi, Pushmeet Kohli, and Alhussein Fawzi. 2024. Mathematical Discoveries from Program Search with Large Language Models. *Nature* 625, 7995 (Jan. 2024), 468–475. doi:10.1038/s41586-023-06924-6

Shivan Jassim. 2023. GRASP: A Novel Benchmark for Evaluating Language GRounding And Situated Physics Understanding in Multimodal Language Models. *arXiv.org* (Nov. 2023). doi:10.48550/arxiv.2311.09048

Haowen Sun, Ruikun Zheng, Haibin Huang, Chongyang Ma, Hui Huang, and Ruizhen Hu. 2024. LGTM: Local-to-Global Text-Driven Human Motion Diffusion Model. In *ACM SIGGRAPH 2024 Conference Papers* (Denver, CO, USA) *(SIGGRAPH '24)*. Association for Computing Machinery, New York, NY, USA, Article 66, 9 pages. doi:10.1145/3641519.3657422

Hao Tang, Darren Key, and Kevin Ellis. 2024. WorldCoder, a Model-Based LLM Agent: Building World Models by Writing Code and Interacting with the Environment. arXiv:2402.12275 [cs.AI] https://arxiv.org/abs/2402.12275

Rodrigo Toro Icarte, Toryn Q. Klassen, Richard Valenzano, and Sheila A. McIlraith. 2018. Reward Machines: Exploiting Reward Function Structure in Reinforcement Learning. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*.

Rodrigo Toro Icarte, Ethan Waldie, Toryn Klassen, Rick Valenzano, Margarita Castro, and Sheila McIlraith. 2019. Learning Reward Machines for Partially Observable Reinforcement Learning. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2019/file/532435c44bec236b471a47a88d63513d-Paper.pdf

Marcell Vazquez-Chanlatte, Susmit Jha, Ashish Tiwari, Mark K Ho, and Sanjit Seshia. 2018. Learning Task Specifications from Demonstrations. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2018/file/74934548253bcab8490ebd74afed7031-Paper.pdf

Lionel Wong, Katherine M. Collins, Lance Ying, Cedegao E. Zhang, Adrian Weller, Tobias Gerstenberg, Timothy O'Donnell, Alexander K. Lew, Jacob D. Andreas, Joshua B. Tenenbaum, and Tyler Brooke-Wilson. 2025. Modeling Open-World Cognition as On-Demand Synthesis of Probabilistic Models. arXiv:2507.12547 [cs.CL] https://arxiv.org/abs/2507.12547

Kun Xiang, Heng Li, Terry Jingchen Zhang, Yinya Huang, Zirong Liu, Peixin Qu, Jixi He, Jiaqi Chen, Yu-Jie Yuan, Jianhua Han, Hang Xu, Hanhui Li, Mrinmaya Sachan, and Xiaodan Liang. 2025a. SeePhys: Does Seeing Help Thinking? - Benchmarking Vision-Based Physics Reasoning. *ArXiv* abs/2505.19099 (2025). https://api.semanticscholar.org/CorpusID:278905129

Kun Xiang, Terry Jingchen Zhang, Yinya Huang, Jixi He, Zirong Liu, Yueling Tang, Ruizhe Zhou, Lijing Luo, Youpeng Wen, Xiuwei Chen, Bingqian Lin, Jianhua Han, Hang Xu, Hanhui Li, Bin Dong, and Xiaodan Liang. 2025b. Aligning Perception, Reasoning, Modeling and Interaction: A Survey on Physical AI. arXiv:2510.04978 [cs.AI] https://arxiv.org/abs/2510.04978

Xinrun Xu, Pi Bu, Ye Wang, Börje F. Karlsson, Ziming Wang, Tengtao Song, Qi Zhu, Jun Song, Zhiming Ding, and Bo Zheng. 2025a. DeepPHY: Benchmarking Agentic VLMs on Physical Reasoning. arXiv:2508.05405 [cs] doi:10.48550/arXiv.2508.05405

Yinggan Xu, Yue Liu, Zhiqi Wu Gao, Changnan Peng, and Di Luo. 2025b. PhySense: Principle-Based Physics Reasoning Benchmarking for Large Language Models. *ArXiv* abs/2505.24823 (2025). https://api.semanticscholar.org/CorpusID:279070750

Haoran Ye, Jiarui Wang, Zhiguang Cao, Federico Berto, Chuanbo Hua, Haeyeon Kim, Jinkyoo Park, and Guojie Song. 2024. ReEvo: Large Language Models as Hyper-Heuristics with Reflective Evolution. (2024). doi:10.48550/ARXIV.2402.01145

Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montse Gonzalez Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, Brian Ichter, Ted Xiao, Peng Xu, Andy Zeng, Tingnan Zhang, Nicolas Heess, Dorsa Sadigh, Jie Tan, Yuval Tassa, and Fei Xia. 2023. Language to Rewards for Robotic Skill Synthesis. arXiv:2306.08647 [cs.RO] https://arxiv.org/abs/2306.08647

Longwen Zhang, Qiwei Qiu, Hongyang Lin, Qixuan Zhang, Cheng Shi, Wei Yang, Ye Shi, Sibei Yang, Lan Xu, and Jingyi Yu. 2023. DreamFace: Progressive Generation of Animatable 3D Faces under Text Guidance. *ACM Trans. Graph.* 42, 4, Article 138 (July 2023), 16 pages. doi:10.1145/3592094

Xinyu Zhang, Yuxuan Dong, Yanrui Wu, Jiaxing Huang, Chengyou Jia, Basura Fernando, Mike Zheng Shou, Lingling Zhang, and Jun Liu. 2025. PhysReason: A Comprehensive Benchmark towards Physics-Based Reasoning. In *Annual Meeting of the Association for Computational Linguistics*. https://api.semanticscholar.org/CorpusID:276422440

Weichao Zhou and Wenchao Li. 2022. Programmatic Reward Design by Example. *Proceedings of the AAAI Conference on Artificial Intelligence* 36, 8 (Jun. 2022), 9233–9241. doi:10.1609/aaai.v36i8.20910

## A  Optimisation Examples

Figure 13 shows example optimisation traces from our system, sourced from the optimisation experiment, demonstrating the effectiveness of our approach.

## B  Learned Library

We separate the learned patterns into two categories: *Guided Patterns* and *Self-discovered Patterns*. Guided patterns were proposed by us to cover a range of common physical interactions we expected to see in the environment. Self-discovered patterns were proposed by the system itself, based on reasoning traces from the Q&A benchmark.

Additionally, in figure 14 we visualise the connectedness of the total learned library i.e. which patterns take other patterns as input.

**Guided Patterns**

**A1.** sliding contact - One object maintains continuous contact with another stationary object while its position changes significantly relative to the other. The stationary object does not move, and the moving object's motion is primarily tangential along the contact surface.

**A2.** rolling over - One object maintains contact with another object while moving a substantial distance relative to it (e.g., above a motion threshold) and then eventually breaks contact. The motion suggests rolling or traversing over the other object rather than just sliding.

**A3.** grid cell transition - The 2D space is discretized into a 10x10 grid, and a pattern is emitted when an object's position crosses from one grid cell to a neighboring grid cell. This captures spatial transitions in the object's trajectory at a coarse spatial resolution.

**A4.** velocity magnitude change - A pattern is emitted when an object's speed changes by a small, medium, or large amount relative to its prior speed. The size of the change is categorized into discrete bins (e.g., small/medium/large) based on predefined thresholds.

**A5.** opposing motion (lever) - One object moves in a given direction causes another object simultaneously to move in the opposite direction, creating a lever effect.

**A6.** moving object hits stationary object - One object is initially at rest while another object moves toward it and collides with it. The collision event is detected when their bodies make contact and the stationary object experiences an impulse or sudden change in velocity.

**A7.** post-collision induced motion - An object that was not directly part of a collision begins to move shortly after the collision occurs, indicating it was supported or constrained by one of the colliding objects. The onset of motion is causally linked in time to the collision but without direct contact during the impact.

**A8.** near collision - Two objects approach each other closely, reaching a separation below a small distance threshold, but never make actual contact. Their relative trajectories suggest an almost-collision that is averted.

**A9.** global rest state - All objects in the simulation are at rest, with velocities below a small threshold for a sustained period. No significant motion or collisions occur during this interval.

**A10.** airborne (no contacts) - An object has zero contacts with any other object (including ground/walls) for a sustained interval, indicating it is fully airborne or unsupported.

**A11.** brief tap - Two objects make contact for only a very short time (below a duration threshold), suggesting a quick tap rather than sustained interaction.

**A12.** bounce / rebound - After colliding, an object reverses its velocity component along the contact normal (sign flip), indicating a bounce.

**A13.** stack formed - One object comes to rest supported on top of another object (stable support contact) for at least a minimum duration.

**A14.** sliding-to-rolling transition - An object in sliding contact transitions to rolling-like motion (consistent angular/linear relationship), indicating it "catches" and begins rolling.

**A15.** tip over - An object that was stationary begins to move towards the ground after a collsion.

**A16.** lose support - An object that was previously supported by another object loses that support, typically when the supporting object moves or is removed.

**A17.** support relationship - Object A is supported by object B: A has contact with B and A's center is above B's center.

**A18.** cell transitions create pattern - A sequence of grid cell transitions occurs for an object, creating a pattern of movement across multiple adjacent cells in the discretized 2D space. x+ is east, y+ is north

**Self-discovered Patterns**

**A1.** first red contact - The first moment after the action when the user-placed red ball initiates contact with any object (first CollisionStart for red).

**A2.** transient goal contact - Green and blue make contact and later separate before the final frame, with TaskComplete.success = False.

**A3.** support handoff beneath object - An object remains in continuous supporting contact while its supporting body changes from one object to another within a short time window, indicating a support handoff.

**A4.** pivoted sweep rotation - An object maintains a nearly fixed contact point (pivot region) with another body and rotates about it by more than a specified angle within a short interval, executing a hinge-like sweep.

**A5.** co-moving bonded pair - Two non-static objects maintain continuous contact for a minimum duration while their relative speed remains below a small threshold and their centers translate significantly, indicating co-moving bonded motion.
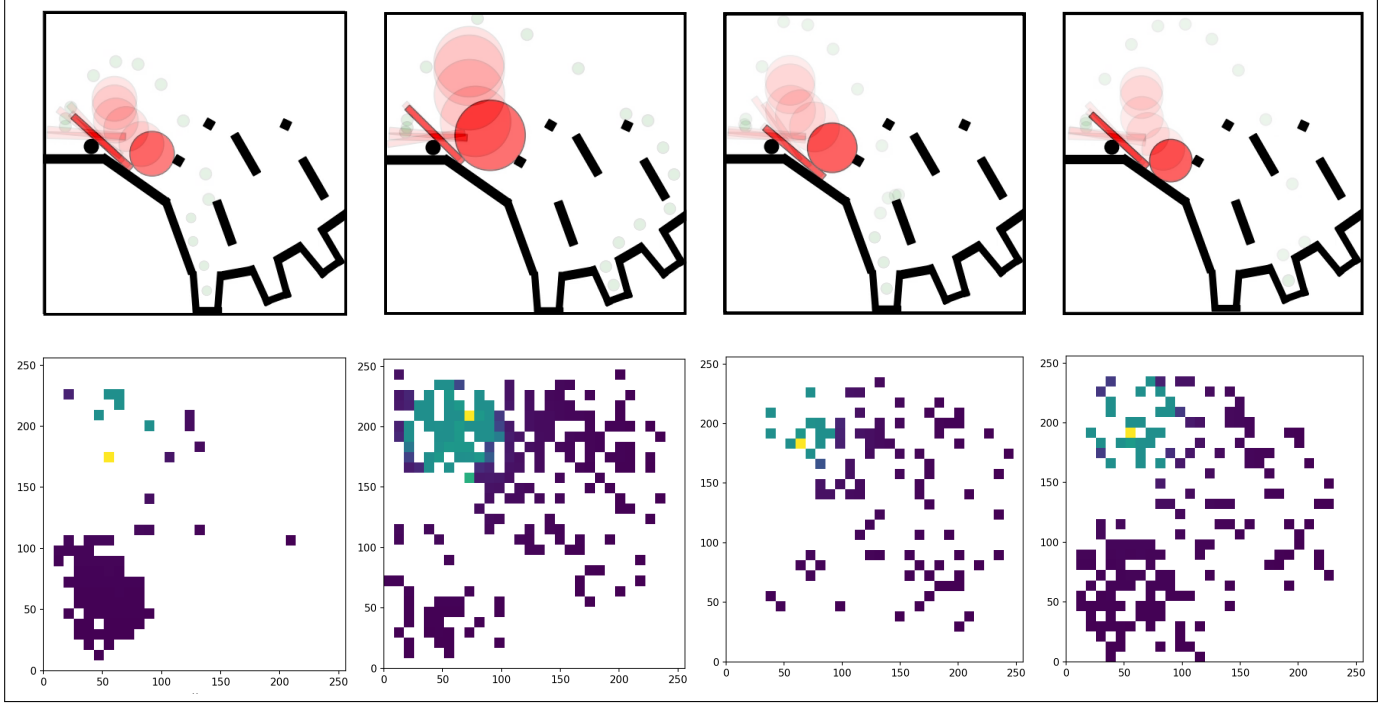
Fig. 13. Example optimisation traces from our system. Each column is a different simulation, where the first row is a visualisation of the final simulation trace, and the second row is a heatmap showing the optimisation search process over 1000 iterations. Brighter colours indicate higher reward program values. In each case the natural language goal was to get the green ball in one of the three buckets. The system successfully synthesised reward programs that guided the optimisation to successful solutions in all cases.

**A6.** goal axis alignment - The centers of the green and blue objects achieve near-equal x or y coordinates (within a small tolerance), indicating a horizontal or vertical axis alignment for the first time in the run.

**A7.** first blue contact - Emit when the first CollisionStart in the run involves the blue object and any other object (ignoring any pre-existing overlap by requiring a prior no-contact frame for blue).

**A8.** last red influence on green - Emit at the last CollisionStart between the red object and the green object in the run, after which no further red–green contacts occur until the simulation ends.

**A9.** red-to-green collision relay - Red collides with an intermediate object X, and within a short time window X collides with green while red never directly contacts green during that window.

**A10.** wedged between two bodies - A dynamic object maintains simultaneous contact with two or more distinct bodies whose average contact normals are roughly opposing, while its speed and positional variance stay below small thresholds for a minimum duration, indicating it is mechanically wedged.

**A11.** high-energy impact - At a CollisionStart, the relative speed between the colliding bodies exceeds a high threshold (e.g., normal component), marking a hard, high-energy impact.

**A12.** goalward acceleration burst - Green's speed increases significantly over a short interval while its distance to blue decreases monotonically over the same interval (beyond set thresholds).

**A13.** sustained goal closing run - For a minimum duration, the green–blue center-to-center distance decreases monotonically while both objects are moving and the instantaneous closing speed stays above a threshold.

**A14.** free-flight goalward approach - Green maintains a continuous no-contact interval during which its center-to-blue distance decreases monotonically by at least a threshold over a minimum duration, optionally culminating in contact.

**A15.** sustained goalward convergence - Over a continuous window, the center-to-center distance between green and blue decreases monotonically by at least a minimum amount and duration, with no increases beyond a small tolerance.

**A16.** side-swap between green and blue - The x-order of the green and blue centers reverses for the first time (green goes from left of blue to right, or vice versa) without requiring contact.

**A17.** first red–blue impact - Emit when the earliest CollisionStart occurs between the red ball and the blue object, regardless of duration or later outcomes.

**A18.** vertical order flip - For a specified pair of objects, the sign of their center y-difference changes, indicating a vertical ordering flip at the frame where their y-coordinates cross within a small tolerance.

**A19.** deflected away from blue - Within a short window after any collision involving green, the angle between green's velocity and the green→blue vector increases by more than a threshold and stays above a misalignment threshold for a minimum duration while their distance stops decreasing.

**A20.** stationary blue success - TaskComplete.success = True occurs while blue's speed stays below a small threshold throughout a window around the first green–blue contact, indicating a stationary-blue goal.

**A21.** moving-blue success - TaskComplete.success = True occurs while blue's speed exceeds a motion threshold within a short window around the first green–blue CollisionStart (using CollisionStart, TaskComplete, and blue's velocity magnitude change or grid cell transition to verify motion).

**A22.** red ricochet bounce - Following a high-energy impact involving red, red reverses its normal velocity component within a brief interval, indicating a ricochet off the impact surface.

**Example Code: "Bounce" Pattern**

```python
def find_pattern(trace):
    """Detect pattern: bounce / rebound
    Return a list of events: {'time': float, 'description': str,
      'parameters': dict}
    Keep timestamps non-decreasing; return [] if not found.
    """
    events = []

    # Collect all collisions that occurred
    collision_events = [(event.time, event.parameters["a_id"],
      event.parameters["b_id"])
                        for event in trace.events
                        if event.uid == "CollisionStart"]

    if not collision_events:
        return events

    # Get object velocities at each frame to detect velocity
      changes
    obj_velocities = {}  # {obj_id: [(time, vel_x, vel_y), ...]}

    for frame in trace.frames:
        for obj in frame.objects:
            if obj.id not in obj_velocities:
                obj_velocities[obj.id] = []
            obj_velocities[obj.id].append((frame.time,
      obj.velocity[0], obj.velocity[1]))

    # Process each collision
    for collision_time, a_id, b_id in collision_events:
        # Get velocity history for both objects around the
          collision time
        a_history = obj_velocities.get(a_id, [])
        b_history = obj_velocities.get(b_id, [])

        if not a_history or not b_history:
            continue

        # Find velocities just before and after collision
          (approximate)
        prev_a = None
        next_a = None
        prev_b = None
        next_b = None

        # Find the frame just before and just after collision
        for time, vel_x, vel_y in a_history:
            if time < collision_time and (prev_a is None or time >
      prev_a[0]):
                prev_a = (time, vel_x, vel_y)
            elif time > collision_time and (next_a is None or time
      < next_a[0]):
                next_a = (time, vel_x, vel_y)

        for time, vel_x, vel_y in b_history:
            if time < collision_time and (prev_b is None or time >
      prev_b[0]):
                prev_b = (time, vel_x, vel_y)
            elif time > collision_time and (next_b is None or time
      < next_b[0]):
                next_b = (time, vel_x, vel_y)

        # Check if a bounced
        if prev_a and next_a:
            prev_vx, prev_vy = prev_a[1], prev_a[2]
            next_vx, next_vy = next_a[1], next_a[2]

            # A bounce is characterized by a significant change in
      velocity direction,
            # especially when one component changes sign
            # We use a simple heuristic: if velocity component
      changes sign and
            # the original magnitude was above threshold, it's
      likely a bounce
            if ((prev_vx * next_vx < 0 and abs(prev_vx) > 1e-6) or
                (prev_vy * next_vy < 0 and abs(prev_vy) > 1e-6)):
                events.append({
                    "time": collision_time,
                    "description": "bounce",
                    "parameters": {
                        "object_id": a_id
                    }
                })

        # Check if b bounced
        if prev_b and next_b:
            prev_vx, prev_vy = prev_b[1], prev_b[2]
            next_vx, next_vy = next_b[1], next_b[2]

            if ((prev_vx * next_vx < 0 and abs(prev_vx) > 1e-6) or
                (prev_vy * next_vy < 0 and abs(prev_vy) > 1e-6)):
                events.append({
                    "time": collision_time,
                    "description": "bounce",
                    "parameters": {
                        "object_id": b_id
                    }
                })

    # Sort by time to ensure non-decreasing order
    events.sort(key=lambda x: x["time"])

    # Remove duplicates - keep only first occurrence of each
      object_id at each timestamp
    seen = set()
    unique_events = []

    for event in events:
        key = (event["time"], event["parameters"]["object_id"])
        if key not in seen:
            seen.add(key)
            unique_events.append(event)

    return unique_events
```

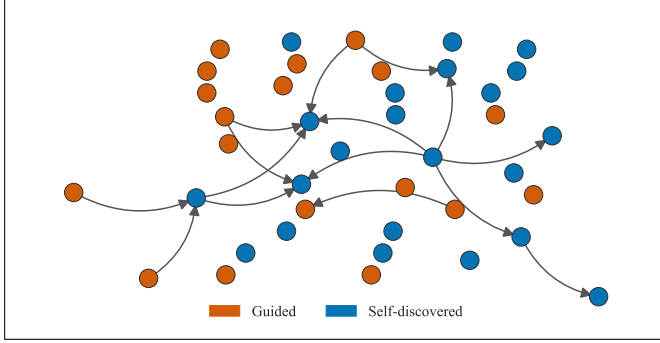Listing 1. Pattern detector: bounce / rebound

Fig. 14. Visualisation of the connectedness of the total learned library i.e. which patterns take other patterns as input.

## C Q&A Task Templates

Below we list the question templates used in our Q&A benchmark. In each case, {PATTERN} is replaced by one of the learned patterns from the library, and {red/green/blue} is replaced by one of the object colours (red, green, or blue) to create a specific question.

**A1.** How many distinct objects does the red/green/blue ball touch between t0 and t1?

**A2.** PATTERN happened with which object(s)?

**A3.** Which object is closest to the red/green/blue ball at the start?

**A4.** Which object blocks the direct line between red/green/blue and red/green/blue (if any)?

**A5.** Which object travels the greatest total distance during the simulation?

**A6.** Which object reaches the highest maximum speed during the simulation?

**A7.** What is the first object that the red/green/blue ball collides with?

**A8.** What percentage of frames contain at least one collision involving the red/green/blue ball (with any object or wall)?

**A9.** What percentage of frames contain a collision between the red/green/blue ball and any other object (excluding walls and ground)?

**A10.** What percentage of frames contain a collision between the red/green/blue ball and the ground?

**A11.** What percentage of frames contain at least one collision between any two moving objects (excluding walls and ground)?

**A12.** What percentage of frames does the red/green/blue ball spend in free fall (not touching any object or wall)?

**A13.** What percentage of frames have at least one object moving in the scene?

**A14.** What percentage of frames have no objects moving in the scene?

**A15.** What percentage of frames do the red/green/blue and red/-green/blue objects touch each other?

**A16.** What percentage of frames is the red/green/blue ball in contact with any wall (left, right, or top boundary)?

**A17.** What percentage of frames is the red/green/blue ball in contact with any static obstacle (excluding ground and walls)?

**A18.** Given the first part of the video, will any collision between a red/green/blue and red/green/blue object occur in the future? Answer yes/no.

**A19.** Given the first part of the video, will the red/green/blue ball ever touch the green object in the future? Answer yes/no.

**A20.** Given the first part of the video, will the red/green/blue ball ever touch the blue object in the future? Answer yes/no.

**A21.** Given the first part of the video, will there be any collision at all between moving objects (excluding walls and ground) in the future? Answer yes/no.

**A22.** Given the first part of the video, will PATTERN happen in the second part? Answer yes/no.

**A23.** Given the first part of the video, will the red/green/blue ball touch the ground for the first time in the future? Answer yes/no.

**A24.** Given the first part of the video, will the red/green/blue ball collide with any wall in the future? Answer yes/no.

**A25.** Given the first part of the video, will the red/green/blue ball cross the vertical line through the green object's initial position in the future? Answer yes/no.

**A26.** Given the first part of the video, will the red/green/blue ball cross the horizontal line through the blue object's initial position in the future? Answer yes/no.

**A27.** Given the first part of the video, will there be any frame in the future in which no objects are moving in the scene? Answer yes/no.

## D Human Survey Details

Figure 15 shows the introduction shown to participants before they began the task. Figure 16 shows an example question shown to participants during the survey. Figure 17 shows an example simulation trace and summaries from the human survey.
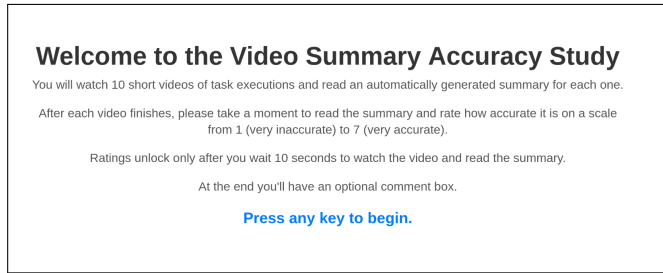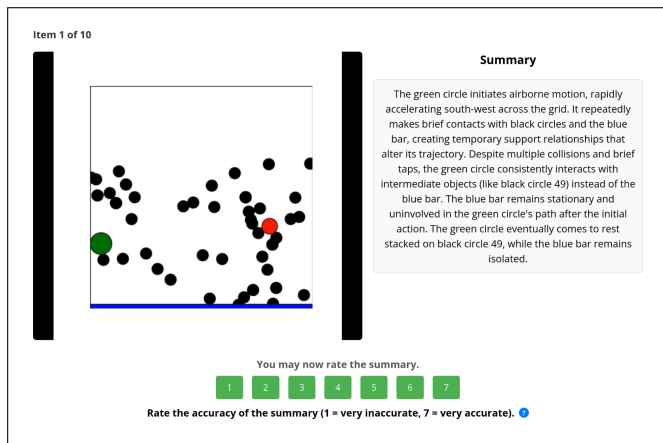
**Welcome to the Video Summary Accuracy Study**

You will watch 10 short videos of task executions and read an automatically generated summary for each one.

After each video finishes, please take a moment to read the summary and rate how accurate it is on a scale from 1 (very inaccurate) to 7 (very accurate).

Ratings unlock only after you wait 10 seconds to watch the video and read the summary.

At the end you'll have an optional comment box.

**Press any key to begin.**

Fig. 15. Survey introduction shown to participants before they began the task.



Fig. 16. Example question shown to participants during the survey.



Fig. 17. Example simulation trace and summaries from the human survey.

Below are examples of the summaries shown for the simulation trace in figure 17.

**Summary A (Annotation Only).**

```
1  The red ball begins moving south-east after a velocity
       change, rolling across the green bar and briefly
       tapping it. This rolling contact propels the green
       bar toward the blue bar. As the green bar moves,
       it eventually makes stable contact with the blue
       bar, forming a persistent stack. The red ball
       continues its path, interacting with the red jar
       and black bars through opposing motion and brief
       taps, but these interactions do not disrupt the
       green-blue contact.
```

**Summary B (Video Only).**

```
1  A red ball is launched toward the scene, colliding with
       the red jar and displacing it. The moving jar then
       strikes the green bar, propelling it laterally.
       The green bar slides across the surface,
       approaching the blue bar. The blue bar adjusts its
       position slightly as the green bar nears,
       maintaining a path toward contact. The green bar's
       continued motion brings its end into direct
       contact with the blue bar's end. This sustained
       interaction ensures stable contact at the final
       frame. The sequence demonstrates a clear causal
       chain: the red ball's impact initiated the jar's
       movement, which transferred energy to the green
       bar, guiding it to meet the blue bar. No obstacles
       interrupted the motion, allowing precise
       alignment. The green and blue bars achieved full
       contact, satisfying the goal with no gap
       remaining. The outcome confirms the simulation
       succeeded due to sequential, energy-transfer
       interactions.
```

**Summary C (Video + Annotation).**

```
1  The red ball begins moving south-east after a velocity
       change. The ball's movement triggers a rolling
       contact with the green bar, briefly supporting it
       before the green bar settles onto the blue bar.
       Subsequent interactions stabilize the green bar
       atop the blue bar, forming a persistent stack.
```

**Survey help guide.**

```
1  An accurate summary should describe the main events in
       the video without fabricating details, match the
       sequence and outcomes shown, and avoid omitting
       key events. In particular, commonly incorrect
       details are: colours, object shapes, sequence of
       events, and positions.
```

## E Language Model Prompts

Below are all of the prompts used throughout our system (besides some minor prompts for error handling and refining outputs).

```
1  Reward DSL reference (parameter-aware):
2
3  - EVENT("uid", {params?}): true if an event with the
       given UID/label occurs. Optional params let you
       require matching event parameters (see library
       parameter schemas below). Example:
       EVENT("abstraction_503681", {"red_ball_id": 8,
       "green_object_id": 7})
4  - AND(expr1, expr2, ...): all child expressions must
       evaluate to true.
5  - OR(expr1, expr2, ...): returns true when at least one
       child expression is true.
```

```
6 - NOT(expr): logical negation.
7 - AFTER("uid_a", "uid_b", min_delta=None,
      max_delta=None, first_params=None,
      second_params=None): true if uid_a occurs after
      uid_b and within optional time bounds; params
      filter each event.
8 - WITHIN("uid_a", "uid_b", window, event_params=None,
      reference_params=None): shorthand for AFTER where
      uid_a occurs no more than `window` seconds after
      uid_b.
9 - COUNT("uid", count, params=None) / GT / LT: count
      occurrences of an event (optionally filtered by
      parameters).
10 - NEARBY_AT(obj_id, x, y, t, threshold_strength=0.1):
      true if object `obj_id` is within
      `threshold_strength * 256` units of point (x, y)
      at simulation time `t` in [0,1]. Use this for
      spatial proximity checks.
11 - OBJECT_ID(color, shape): returns the object ID for
      the object with the given color and shape in the
      scene (e.g., OBJECT_ID("red", "circle")). Use this
      to get IDs for NEARBY_AT or EVENT parameters.
      Shape may be "bar", "circle", "jar",
      "standingsticks", or "any"; color may be "red",
      "blue", "green", or "black".
12
13 Examples:
14     AND(
15         EVENT("abstraction_217640"),
16         AFTER("abstraction_217640",
      "abstraction_307800", first_params={"red_ball_id":
      8}, second_params={"green_object_id": 7}),
17         NOT(EVENT("abstraction_612355", {"frame_index":
      44})))
18     )
19     OR(
20         COUNT("abstraction_612355", 2, {"red_ball_id":
      3}),
21         GT("abstraction_661256", 3)
22     )
23     NEARBY_AT(OBJECT_ID("red", "circle"), 100.0, 150.0,
      0.5, threshold_strength=0.2)
24
25 IMPORTANT:
26 - Identifiers are case-insensitive when matching event
      labels; UIDs are matched exactly.
27 - Parameter matching requires ALL provided keys to
      match the emitted event parameters (strings are
      case-insensitive; numbers match by value). Leave
      params empty to match any.
28 - Identifiers, such as NEARBY_AT, should not have their
      arguments included in the expresseion, i.e.
      correct use is NEARBY_AT(5, 100.0, 150.0, 0.5),
      not NEARBY_AT(obj_id=5, x=100.0, y=150.0, t=0.5).
29 - Identifiers MUST exist in the library; DO NOT invent
      new ones.
30 - In addition to the library events, there are three
      important built-in events:
31     - "uid": "CollisionStart", emitted when two objects
      start colliding; parameters include "a_id" and
      "b_id".
32     - "uid": "CollisionEnd", emitted when two objects
      stop colliding; parameters include "a_id" and
      "b_id".
33     - "uid": "TaskComplete", emitted when the
      simulation ends, successful if the green and blue
      objects are touching; parameters include "success"
      (boolean).
34 - See the event library documentation for parameter
      schemas for built-in and abstraction events.
```

Listing 2. Domain Specific Language Specification for LM

```
1 Given the following goal and the DSL reference, propose
      one DSL expression. This expression defines a
      reward function. This reward function will be
      maximised by an optimisation process. Do not
      include any comments in the DSL code.
2 Events can expose parameters; use them to precisely
      target entities (see library parameter schemas
      below).
3
4 TIPS:
5 - It is easier for the optimisation process to maximise
      the reward function when it is expressed as a sum
      of positive terms i.e. a step by step addition of
      clauses with AND operators.
6 - Think about how you want to achieve the goal, and
      express that in the reward function. Don't just
      describe a single end condition, but build up the
      reward function step by step.
7 - Try to reason about the video that you see, this is
      the state that the optimisation process will be
      acting in.
8 - DO NOT include comments in the DSL code.
9
10 Goal:
11 {{ goal }}
12
13 DSL reference:
14 {{ dsl_guide }}
15
16 Library summary:
17 {{ library_summary }}
18
19 Scene summary:
20 {{ scene_summary }}
21
22 Respond in think/answer blocks, with ONLY the DSL
      inside ```dsl fences.
23 <think>reason carefully referencing evidence</think>
24 <answer>```dsl
25 ... your DSL ...
26 ```</answer>
```

Listing 3. LM Reward Program Synthesis Prompt

```
1 You will improve a Python function find_pattern(trace)
      that scans a physics Trace and returns a list of
      events.
2 Each event MUST be a dict: {"time": float,
      "description": str, "parameters": dict}
3
4 Pattern label (stable identifier): {{ label }}
5 Event description (code should detect this pattern):
6 {{ description }}
7
8 {{ trace_spec }}
9
10 Existing library formattings (i.e. the event paramaters
      in event objects in trace.events):
11 {{ formattings }}
12 These events may occur in the trace.events list. If
      they are used in the code, make sure to explicitly
      reference the uid of the event in the code.
13
14 Constraints:
15 - Implement: def find_pattern(trace): -> list[dict]
16 - Each dict of the returned list should have keys:
      "time", "description", "parameters"
17 - No imports, no I/O, no eval/exec
18 - Use math.<fn> if needed (assume math provided)
19 - Prefer simple loops/thresholds; keep code focused and
      efficient
20 - Return [] if nothing is detected
21 - Make sure the code does not detect a pattern every
      frame, this is always incorrect behaviour
22 {{ extra_constraints }}
23
```

```
24 It is very important that your detector focuses on
       detecting NEW patterns that are not already
       covered by existing library abstractions i.e. do
       not simply reuse existing events but perform some
       new logic to detect novel patterns while utilizing
       existing events if helpful.
25
26 Current code:
27 ```python
28 {{ parent_code }}
29 ```
30
31 Errors / issues with current code:
32 {{ errors }}
33
34 After the code block, output a JSON description of the
       expected "parameters" shape returned by the
       detector.
35 Each key in the JSON should be a parameter name and its
       value must be a string describing the type (e.g.
       "tuple[float, float]").
36 Return exactly three outputs in this order:
37 1. Let's think step by step...
38 2. ```python ... ```
39 3. ```json ... ```
```

Listing 4. LM Code Evolution Prompt

```
1 You are generating reusable event pattern labels and
       descriptions for a 2D physics simulation. These
       patterns will help analysts understand key moments
       in simulation traces.
2
3 Simulation domain (information given to analysts):
4 {{ trace_spec }}
5
6 Objective:
7 - The ONLY goal is for the green object to touch the
       blue object.
8
9 What patterns are:
10 - Named, reusable events that help a human visualize
       key moments in a trace.
11 - Code will be written later to detect these patterns;
       your output seeds those detectors.
12
13 Current library snapshot (UID, label if any,
       description):
14 {{ library_table }}
15
16 RL reasoning (<think> snippets that reveal what
       analysts care about):
17 {{ rl_thinks }}
18
19 Your task:
20 {{ abstract_guidance }}
21 - Propose {{ K }} NEW patterns (not duplicates or
       near-duplicates of existing ones).
22 - Make sure each pattern is distinct and captures a
       unique aspect of the simulation traces.
23 - It is very important that suggested patterns are not
       similar to existing ones in the library.
24 - Each pattern:
25   - reason: why this pattern is useful given the
       objective and the reasoning
26   - description: one sentence describing the pattern in
       detail
27   - label: a short, descriptive phrase (3-7 words),
       importantly it should be scene-agnostic
28
29 Output:
30 Make sure to think about what patterns should be
       suggested first, then output a JSON array like
       this:
31 ```json
32 [
```

```
33   {"reason": "one or two sentences", "description":
       "one sentence", "label": "short sentence"},
34   ...
35 ]
36 ```
```

Listing 5. LM Label Suggestion Prompt

```
1 This domain consists of various objects in a 2D physics
       environment. A user takes an action by placing a
       red ball in the scene, with a given (x, y)
       position and radius. The goal of the simulation is
       to have the object coloured green touch the object
       coloured blue at the end of the simulation. This
       is the only objective. The green and blue object
       can be of any shape (circle, bar, jar,
       standingsticks) and located anywhere in the scene.
       The bounds of the scene are 256x256 units.
2
3 The trace consists of a sequence of frames, each
       representing the state of the environment at a
       specific time step. Each frame contains
       information about the positions, velocities, and
       other properties of all objects in the scene. The
       trace also includes events that occur during the
       simulation, such as collisions between objects or
       when the task is completed successfully.
4
5 All traces have a list of events containing
       'CollisionStart' and 'CollisionEnd' events, which
       indicate when two objects start and stop
       colliding, respectively. The final event in any
       trace is always a 'TaskComplete' event, which
       contains ['parameters']['success'] set to True if
       the task was completed successfully (i.e., the
       green object touched the blue object) or False
       otherwise.
6
7 Trace:
8     action: Action
9     scene: Scene
10    frames: List[Frame]
11    events: List[Event]
12
13 Action:
14    position: List[float] # (x, y) position where the
       red ball is added
15    radius: float         # radius of the red ball
16
17 Scene:
18    objects: List[Object] # Objects present in the
       scene before the simulation starts, both dynamic
       and static objects
19
20 Frame:
21    time: float           # A time between 0.0 and 1.0
       indicating when this frame occurs in the
       simulation, where 0 is the first frame and 1 is
       the final frame
22    objects: list[Object] # Objects present in the
       scene at this time step, dynamic objects only
23
24 Object:
25    description: str      # e.g., "green-circle-2",
       "blue-bar-1", etc.
26    id: int               # unique identifier for the
       object, shown in description, e.g. 2 in
       "green-circle-2"
27    type: str             # one of "circle", "bar",
       "jar", "standingsticks"
28    color: str            # one of "green", "blue",
       "red", "black"
29    velocity: list[float] # (x, y) velocity of the
       object at this frame
30    angle: float          # angle of the object at
       this frame as a float in radians
```

```
31    static: bool          # whether the object is
      static (True) or dynamic (False), when static the
      object does not move and is only listed in
      trace.scene.objects and not in
      trace.frames[:].objects
32    obj_params: dict      # parameters defining the
      object's shape
33
34    if type == "circle":
35        obj_params:
36            center: list[float] # (x,y) position of
      circle center
37            radius: float
38    elif type in ["bar", "jar", "standingsticks"]:
39        obj_params:
40            polygon_0: list[list[float]] # list of
      (x,y) vertices (note there is no 'center' or
      'radius' for these shapes)
41            polygon_1: list[list[float]]
42            ... etc.
43
44    Example:
45    obj_params = {
46        "center": [100.0, 150.0],
47        "radius": 10.0
48    }
49    or
50    obj_params = {
51        "polygon_0": [[x1, y1], [x2, y2], ...],
52        "polygon_1": [[x1, y1], [x2, y2], ...],
53    }
54
55 Event:
56    time: float         # A time between 0.0 and 1.0
      indicating when the event occurred in the
      simulation, where 0 is the first frame and 1 is
      the final frame
57    uid: str            # Unique identifier for the
      event type, e.g., "CollisionStart",
      "CollisionEnd", "TaskComplete", or a library
      abstraction uid like "abstraction_1", etc.
58    parameters: dict    # Parameters associated with the
      event, varying based on event type
59
60    if uid in ["CollisionStart","CollisionEnd"]::
61        parameters:
62            a_id: int # object a id
63            b_id: int # object b id
64            contact_points: list[list[float]]
65    elif uid == "TaskComplete":
66        parameters:
67            success: bool
68    else:
69        # Other abstraction events may have different
      parameters depending on the abstraction
      definition, see library for details
70
71 Examples of API:
72
73    scene_objects = trace.scene.objects
74    for obj in [o for o in scene_objects if not
      o.static]:
75        for frame in trace.frames:
76            objs_in_frame = [o for o in frame.objects
      if o.id == obj.id]
77            ...
78
79    events = trace.events
80    for event in events:
81        uid, time, params = event.uid, event.time,
      event.parameters
82        ...
83
84    if event.uid == "CollisionStart":
85        a_id = event.parameters["a_id"]
```

```
86        obj = [o for o in trace.scene.objects if o.id
      == a_id][0]
87        if obj.color == "green" and obj.type ==
      "circle":
88            radius = obj.obj_params["radius"]
89            ...
```

Listing 6.  Simulation Trace Specification Prompt

## F  Code Evolution Details

**FunSearch Algorithm**

Algorithm 3 outlines the FunSearch procedure used for program synthesis via LLMs. The algorithm maintains multiple islands of program candidates, periodically resetting lower-performing islands to promote diversity and exploration. We adapted this method by providing our own Evaluation function tailored to our code-learning task.

---

**ALGORITHM 3:** FunSearch

---

**Input:** evaluation function Evaluate$(\cdot)$,
  initial (or skeleton) program $g_0(\cdot)$,
  LLM$(\cdot)$,
  number of islands $I$, prompt size $s$, reset period $T_r$
**Output:** Best program found $g^\star(\cdot)$
**for** $i \leftarrow 1$ **to** $I$ **do**
  $\mathcal{D}_i \leftarrow \{g_0\}$
Initialize program database as islands $\mathcal{D} \leftarrow \{\mathcal{D}_i\}_{i=1}^I$;
**for** *iteration* $t \leftarrow 1$ **to** *budget* **do**
  Sample an island $\mathcal{D}_i$ (favor islands with higher best score);
  Sample $k$ best programs $g_1, \ldots, g_s$ from $\mathcal{D}_i$ ;
  $g_{new} \leftarrow$ LLM(BuildPrompt$(g_1, \ldots, g_s)$);
  **if** $g_{new}$ *is valid* **then**
    $v \leftarrow$ Evaluate$(g_{new})$;
    Add $(g_{new}, v)$ to island $\mathcal{D}_i$;
  **if** $t \bmod T_r = 0$ **then**
    Identify worst half of islands;
    Reinitialize lower half of islands by cloning a top program ;
Return $(g^\star, v^\star) \leftarrow \arg\max_{(g,v)\in\mathcal{D}} v$;

---

**ComputeLengthPenalty**

We apply a logarithmic penalty based on code length. We count the total number of lines and compute the penalty as,

$$\lambda = \log(\texttt{num\_lines})/5.$$

To avoid unbounded growth, we cap `num_lines` at 1000, which yields a maximum penalty of approximately 1.5.

**ComputeTimePenalty**

We use a time-based penalty derived from the average annotation time of existing patterns in the library. Let $t$ be the average annotation time for the new pattern and let $\mu$ be the mean annotation time across existing patterns. If $t \leq \mu$, the penalty is 0. If $t > \mu$, we apply a linear penalty that increases from 0 to 1 as $t$ rises from $\mu$ to $2\mu$; specifically, the penalty reaches 1 when the new pattern takes twice the mean time. This value is the maximum penalty—any slower pattern (i.e., $t \geq 2\mu$) receives a penalty of 1.

## G Value Model Goals

Below are the full set of goals used for value model training in our experiments. The goals for hand written rewards are a subset of those used for reward program synthesis.

**Goals for Reward Program Synthesis**

**A1.** Get the green and blue objects to be touching at the end of the simulation. Make the green object touch the blue object, while ensuring the red ball never touches the green object.

**A2.** Make the green object touch the blue object, but make sure the red ball does not touch the green object within the first 20% of the simulation.

**A3.** Make the red ball touch the green object after the green object has already touched the left wall, while preventing the green object from touching the right wall.

**A4.** Get the green object to touch both the blue object and the red ball (in any order), while ensuring the red ball never touches the blue object.

**A5.** Make the green object and blue object touch exactly once, and require that the red ball touches neither of them until after that single touch has occurred.

**A6.** Create a 'near-miss' challenge: around time 0.5 the red ball must be very close to the blue object, but the red ball must never actually touch the blue object at any point.

**A7.** The red ball must touch the green object and the blue object the same number of times (at least once each).

**A8.** Make the green object touch the blue object only after the red ball has contacted the ground and stopped moving.

**A9.** Force the red ball to knock the green object into the blue object exactly once, with no other contacts involving red.

**A10.** Make the green object touch the blue object while preventing any object from contacting the scene boundaries (left, right, or top walls) at any time during the simulation.

**Goals for Hand Written Rewards**

**A1.** Make the green object and blue object touch exactly once, and require that the red ball touches neither of them until after that single touch has occurred.

**A2.** The red ball must touch the green object and the blue object the same number of times (at least once each)

**A3.** Make the green object touch the blue object only after the red ball has contacted the ground and stopped moving.

**A4.** Force the red ball to knock the green object into the blue object exactly once, with no other contacts involving red.

**A5.** Make the red ball touch the green object after the green object has already touched the left wall, while preventing the green object from touching the right wall.