

EXTENDING TEST-TIME SCALING: A 3D PERSPECTIVE WITH CONTEXT, BATCH, AND TURN

Anonymous authors

Paper under double-blind review

ABSTRACT

Reasoning reinforcement learning (RL) has recently revealed a new scaling effect: test-time scaling. Thinking models such as R1 and o1 improve their reasoning accuracy at test time as the length of the reasoning *context* increases. However, this effect is fundamentally limited by the context window of the base models, which remains orders of magnitude smaller than the amount of tokens consumed during training. We revisit test-time enhancement techniques through the lens of scaling effect and introduce a unified framework of multi-dimensional test-time scaling to *extend* the capacity of test-time reasoning. Beyond conventional context-length scaling, we consider two additional dimensions: *batch scaling*, where accuracy improves with parallel sampling, and *turn scaling*, where iterative self-refinement enhances reasoning quality. Building on this perspective, we propose 3D test-time scaling, which integrates context, batch, and turn scaling. We show that: (1) each dimension demonstrates a test-time scaling effect, but with a bounded capacity; (2) combining all three dimensions substantially improves the reasoning performance of challenging testbeds, such as IOI, IMO, and CPHO, and further benefits from human preference feedback; and (3) the human-in-the-loop framework naturally extends to a more open-ended domain, i.e., embodied learning, which enables the design of humanoid control behaviors.

1 INTRODUCTION

Recent progress in reasoning reinforcement learning has introduced a new form of scaling effect by training thinking models such as R1 (Guo et al., 2025) and o1 (OpenAI, 2024). Unlike conventional models that directly map input to output, a thinking model performs intermediate reasoning computation before producing its final answer. A striking phenomenon emerges during the reinforcement learning process: as the model is trained to reason over progressively longer contexts, its reasoning accuracy steadily improves (Shi et al., 2025; Aggarwal & Welleck, 2025). At test time, this trend continues: extending the reasoning context length consistently leads to higher accuracy. This phenomenon is referred to as test-time scaling of reasoning models (Muennighoff et al., 2025).

However, the potential of test-time scaling is fundamentally constrained by the context window size of current models. Even the most advanced commercial reasoning systems today support fewer than one million tokens of context—negligible compared with the scale of training-time compute, where tens of trillions of tokens are typically consumed during pretraining or post-training. This discrepancy naturally raises a question:

How should we extend the capacity of test-time scaling?

Notably, there have been many popular heuristics to enhance the reasoning model’s performance at test time. For example, majority voting improves accuracy by generating multiple candidate outputs in parallel and selecting the most frequent one (Wang et al., 2023a). Other approaches, such as Reflexion (Shinn et al., 2023) and in-context learning (Madaan et al., 2023), perform iterative self-refinement, where a model repeatedly revisits and improves its own solutions. Empirically, taking multiple refinement steps leads to a higher accuracy compared with directly outputting the solution.

In this paper, we revisit these diverse techniques within a unified framework of *multi-dimensional test-time scaling*. Specifically, we consider three dimensions: (1) Context scaling: reasoning accuracy improves with longer thinking context lengths; (2) Batch scaling: methods such as majority

054 vote can be viewed as scaling along a batch dimension, where more parallel samples yield better
055 aggregated answers; (3) Turn scaling: iterative refinement methods correspond to scaling along a
056 turn dimension, where more refinement turns enhance accuracy. Each of these dimensions of scal-
057 ing interacts with the context-length limits and capabilities of base LLMs, creating unique empirical
058 trade-offs.

059 Building on this perspective, we propose *3D test-time scaling*, which integrates all three dimensions:
060 context, batch, and turn. We demonstrate that this unified view substantially extends the ceiling of
061 test-time scaling compute and further enables a human-in-the-loop framework that applies to even
062 open-ended domains.

- 064 • We establish that each scaling dimension individually exhibits a test-time scaling effect:
065 higher token consumption leads to higher accuracy. However, clear scaling limits can be
066 observed for each dimension.
- 068 • We show that the unified 3D test-time scaling is capable of leveraging substantially more
069 tokens for improved reasoning and achieving gold-level performances on challenging
070 Olympiad competition problems, such as IOI, IMO, and CPHO. The framework also ex-
071 tends to a human-in-the-loop setting, where a human operates along the batch dimension
072 and selects the best candidate to further amplify final accuracy.
- 074 • Finally, we extend this human-in-the-loop framework to embodied learning, demonstrating
075 that multi-dimensional test-time scaling enables models to interactively design open-ended
076 behaviors in humanoid robot control.

078 2 RELATED WORK

081 **Scaling Laws.** The cross-entropy loss in large language model pretraining has been shown to
082 scale predictably with key training resources, including model size, dataset size, and compute budget
083 (Kaplan et al., 2020; Rae et al., 2022; Hoffmann et al., 2022). With the emergence of thinking models
084 such as DeepSeek-R1 (Guo et al., 2025) and OpenAI o1 (OpenAI, 2024), researchers investigated
085 scaling laws beyond the number of training tokens. For example, Shi et al. (2025) examines scaling
086 behavior with respect to context length. Scaling laws have also been studied at test-time: Wu et al.
087 (2025); Snell et al. (2024) analyze how performance scales with respect to inference compute under
088 different inference strategies such as majority voting and tree search, as well as tradeoffs between
089 model size and test-time token budgets. In this paper, we focus on test-time scaling and propose
090 a unified framework for characterizing its effects across three dimensions: context scaling, batch
091 scaling, and turn scaling. In contrast, prior work on test-time scaling laws has typically examined
092 only a subset of these aspects.

094 **Test-Time Scaling.** Test-Time Scaling (TTS) refers to the class of algorithms for improving the
095 model’s performance through scaling inference-time compute. TTS methods can be broadly categor-
096 ized into three approaches. *Context scaling* methods improve performance through longer output
097 sequences, exemplified by Chain-of-Thought prompting (Wei et al., 2023), which elicits step-by-
098 step reasoning in large language models to improve performance on various benchmarks. Recent
099 advances in reasoning models like o1 (OpenAI, 2024) and DeepSeek-R1 (Guo et al., 2025) further
100 incentivize this ability, highlighting context scaling as an effective strategy for improving test-time
101 performance. *Batch scaling* approaches leverage parallel computation to explore multiple reason-
102 ing paths. Majority voting is a representative technique that leverages the power of parallel sam-
103 pling (Wang et al., 2023a) by generating multiple independent reasoning paths and selecting the ma-
104 jority final answer. Other work further incorporates test-time search (Yao et al., 2023), Monte-Carlo
105 tree search (Zhang et al., 2024; Xie et al., 2024), and parallel thinking (Ning et al., 2023) to improve
106 the performance. *Turn scaling* methods improve performance through iterative refinement, includ-
107 ing Self-Refine (Madaan et al., 2023), which enables models to iteratively improve outputs through
self-feedback without additional training, and Reflexion (Shinn et al., 2023), which reinforces lan-
guage agents through linguistic feedback and episodic memory to enhance future decision-making.

3 FORMULATION OF TEST-TIME SCALING

LLM Reasoning. In this work, we focus on LLM reasoning. Given a question x , the goal is to derive a correct step-by-step solution y . We assume the existence of a ground-truth verifier \mathcal{R} that evaluates the correctness of a solution y . This verifier \mathcal{R} could have different implementations depending on the specific case in practice. For example, in mathematical reasoning tasks where the goal is to derive a single numerical answer, this verifier could return a 0-1 score indicating whether the answer in solutions y matches the ground-truth answer. In coding tasks, the score is determined by the set of tests passed by the submitted code in solution y . An LLM π_θ is a policy parameterized by θ . Given an input question x , the LLM auto-regressively generates an array of tokens one by one. For a dataset of questions \mathcal{D} , the expected score of an LLM policy π_θ given a question x is defined as $J(\mathcal{D}, \pi_\theta) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)}[\mathcal{R}(x, y)]$.

Test-Time Scaling. Test-time scaling approaches aim to achieve a better score through spending more test-time compute. For instance, context scaling allows the LLM to generate longer responses to conduct in-depth exploration. The efficacy of any test-time scaling method must be evaluated along two key aspects: the expected score and the computational cost. In this work, we quantify computational cost using the theoretical maximum number of tokens generated throughout the inference process.

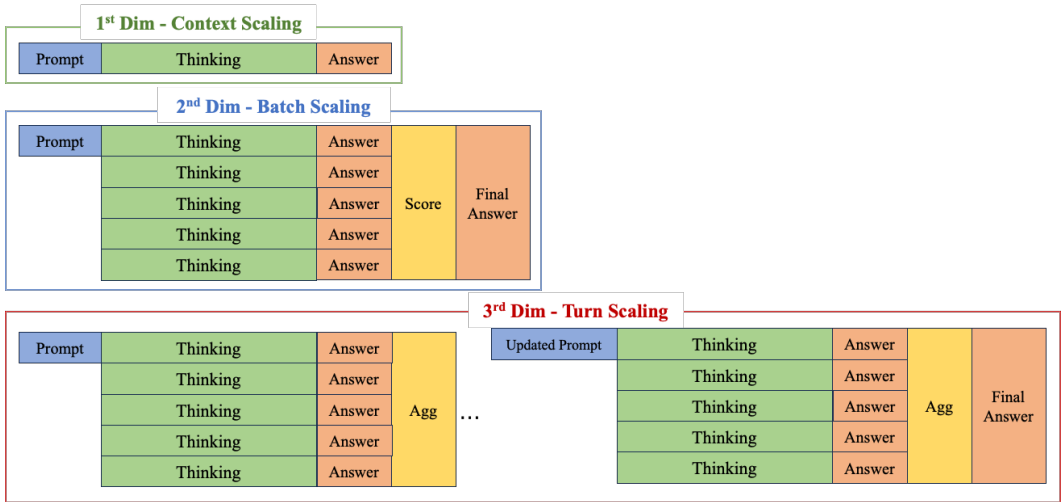


Figure 1: Illustration of Test-time Scaling across three dimensions: context, batch, and turn.

3.1 TEST-TIME SCALING WITH CONTEXT, BATCH, AND TURN

Context Scaling. We first consider the context dimension. In context scaling, we follow the most straightforward approach of controlling the reasoning process with a maximum token budget C . Specifically, when the LLM generates a response $y \sim \pi_\theta(\cdot|x)$ exceeding C , we directly truncate the response. Therefore, the expected score of context scaling under a context length C is $J_{\text{context}}(\mathcal{D}, \pi_\theta, C) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(\cdot|x)}[\mathcal{R}(x, y:C)]$

Batch Scaling. To scale test-time compute along the batch dimension, notable examples of batch scaling approaches include Majority Voting and Best-of-N (Wang et al., 2023b; Snell et al., 2024; Cobbe et al., 2021). Given the size of the batch dimension B , i.e., the number of responses generated in parallel, and the context length C , B responses y_1, y_2, \dots, y_B are first generated under a token budget of C . Then, a final solution is derived by selecting the best response through a scoring function $\text{Score}(y)$, i.e. $y_{\text{final}} = \arg \max_i \text{Score}(y_i)$. In cases where the goal is to derive a single numerical answer, the scoring function could be implemented by counting the occurrence of each candidate answer in the B responses. In coding tasks or mathematical proof tasks where it is infeasible to apply voting, we employ the LLM π_θ to select a single solution out of the B responses. In this paper, we refer to majority voting as **Batch Scaling (Vote)** and the Best-of- N strategy as **Batch Scaling (Best-of- N)**.

3D Scaling with a Turn Dimension. We introduce an additional turn axis to batch scaling, leading to a unified 3D scaling framework, as shown in Fig. 1. In 3D scaling, the whole process takes T turns. Each turn t starts from prompt p_t , containing both the original problem and past experiences, and generates B independent responses within a context length of C , $y_i^t \sim \pi_\theta(\cdot|p_t)$. Similar to the scoring function in batch scaling, an aggregation function $\text{Agg}(y_1^t, \dots, y_B^t)$ is used to aggregate the B responses and generate a *context summary* that serves as a compact compression of the generation history in the first t turns. The prompt for the next turn, p_{t+1} , is then composed by concatenating the input question x and the context summary. The final answer is extracted from the aggregated result of the last turn T .

3D Scaling Variants. Different choices of aggregation functions and 3D configurations result in different algorithmic realizations of 3D scaling. We consider three major variants.

- **Turn Scaling (Reflection):** This approach is a special case of 3D scaling where only one response is generated in each turn, i.e. $B = 1$. In this case, the aggregation function only takes one response as input. We use an LLM to generate a reflection on the response. The reflection would be used as additional feedback for future turns.
- **3D Scaling (LLM Judge):** In this approach, a batch of $B > 1$ responses are generated in each turn. To automatically identify the best response within a batch, we adopt a straightforward method: we input all responses to the Gemini LLM and instruct it to return the index of the optimal one. The responses not selected as the best are then randomly sampled to serve as negative examples.
- **3D Scaling (Human Judge):** Finally, we also consider an interesting instantiation of 3D scaling in a human-in-the-loop manner. In this approach, we use feedback from human experts as the aggregation function, where an expert evaluates the batch of responses and selects the best path forward. This approach is effective in cases when the LLM can not provide an accurate judgment to select the most salient response.

We also remark that it is also feasible to query the LLM to generate complex feedback for future turns, such as summarizations and reflections over the batch (Shinn et al., 2023; Huang & Yang, 2025). For simplicity, in this work, we choose selection as the aggregation function in each turn.

4 EXPERIMENTS

We begin with the experiment setup and then proceed with three evaluation stages. First, we examine the performance of different test-time compute configurations over three dimensions on IMO problems to illustrate the test-time scaling phenomena. Next, we explore how the unified 3D scaling pushes the reasoning capacity on a collection of challenging Olympiad problems. Finally, we extend the framework to a more open-ended setting, embodied learning. With human feedback in the loop, 3D scaling produces robotic control behaviors that are more aligned with human preferences.

4.1 EXPERIMENT SETUP

Base Reasoning Model: We conduct all experiments using Gemini 2.5 Pro Comanici et al. (2025) as the backbone model, chosen for its strong reasoning and coding capabilities in complex problem-solving tasks. To ensure reproducibility, the temperature is fixed at 0.1, yielding highly deterministic outputs across trials. For each domain, we further design tailored system prompts for solution generation and feedback learning; full prompt details are provided in Appendix D.2.

Testbeds: We explore the scaling effect on two types of testbeds:

- **Reasoning Problem-Solving Tasks:** This testbed focuses on rigorous reasoning and algorithmic problem-solving. (1) *Math and Physics Olympics:* We adopt problems from the IMO International Mathematical Olympiad (2025) and CPHO Chinese Physics Olympiad (2022) to evaluate the LLM’s reasoning capabilities. (2) *Coding:* IOI 2025 problems International Olympiad in Informatics (2025) are used to assess programming ability under 3D Scaling. Unlike human contestants who receive submission feedback, the LLM must directly solve tasks without intermediate guidance.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

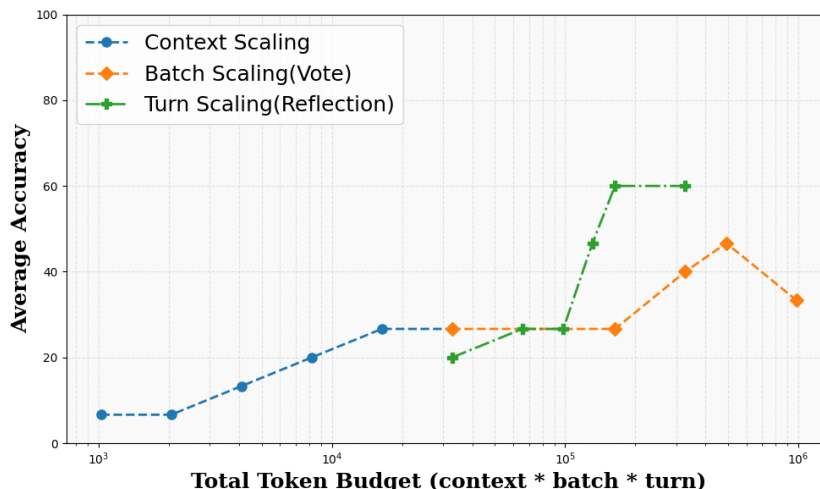


Figure 2: The average accuracy over the IMO2025 dataset as a function of the total thinking budget for individual scaling on three dimensions: context, batch and turn. All three scaling methods achieve substantial improvements at small scales but saturate as the scale becomes larger.

- Innovative Tasks:** This testbed targets embodied AI and emergent behaviors. We use several robotics reinforcement learning tasks from GPU-based IsaacGym Makoviychuk et al. (2021) that cover diverse environments. We also introduce a new task, *HumanoidJump*, which aims to make a humanoid jump in a human-like manner. Designing a reward for this task is an open challenge because human-like jumping lacks easily quantifiable criteria.

Evaluation: For IMO and CPHO problems, every LLM-generated solution is rigorously verified by *human experts* following the scoring guidance. A solution is considered correct only if both the final answer and the entire reasoning process are mathematically valid. For IOI problems, the score is measured over the official IOI test cases. For innovative tasks, we recruit human volunteers to vote for their preferred behaviors.

Human-in-the-loop Feedback: In the setting of 3D scaling with $B > 1$, in addition to using an LLM judge, we can also introduce a human judge to select the best solution among all parallel candidates in each refinement turn according to the task objective. Details about evaluators are presented in Appendix A.4.

4.2 MULTI-DIMENSIONAL TEST-TIME SCALING ANALYSIS

In this subsection, we study test-time scaling on the IMO benchmark. We select three moderately difficult problems (1, 3, and 5), excluding those that are too easy or too hard. Each is tested over five trials, and we report the average accuracy over 3 problems. To fully exploit the backbone LLM, all experiments except *context scaling* fix the context length at 32K.

4.2.1 SINGLE-DIMENSION SCALING ANALYSIS

We investigate the test accuracy by scaling along each of the three dimensions. For **Context Scaling**, we vary context length C from 1k to 32k. For **Batch Scaling (Vote)**, we take the 32K context length with B parallel rollouts ranging from 1 to 30. For **Turn Scaling (Reflection)**, we adopt full context and $B = 1$ while allowing the model to take 1 to 10 refinement turns. Fig. 2 reports the average accuracy as a function of total thinking budget across the three individual scaling dimensions. Performance improves at small scales but quickly plateaus, with little or no gain from further scaling. Notably, batch scaling even degrades at $B = 30$. Note that in IMO problems, the solution is considered correct only if the answer and the process are both correct. Therefore, we hypothesize that this is because majority voting amplifies systematic biases: when the model favors a specific wrong pattern in the derivation process, more samples may reinforce the process bias. We also report this interesting finding as an open question to the community.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

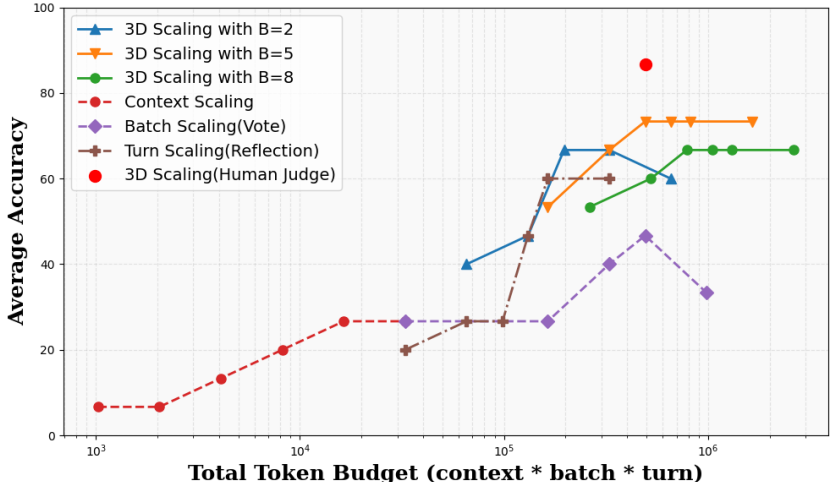


Figure 3: The average accuracy over the IMO2025 dataset as a function of the total thinking budget for individual scaling and 3D Scaling with different batch sizes. 3D Scaling achieves performance beyond the limits of individual scaling, reaching 73.3%. The red marker denotes 3D Scaling with a human judge, which attains 86.7% accuracy, highlighting the effectiveness of human feedback.

4.2.2 3D SCALING ANALYSIS

We conducted 3D Scaling experiments that combine batch scaling and turn-based scaling, using a simple preference aggregation function provided by the LLM Judge. The four solid curves in Fig. 3 about 3D Scaling with various batch sizes show how model accuracy varies with two parameters: the *batch size* B and the number of *turns* T .

The results largely align with observations from single-dimension scaling. Increasing the number of turns T initially improves performance by stabilizing predictions; however, further increases lead to saturation and may even reduce accuracy, likely because an incorrect judgment in one turn by the LLM can adversely affect refinements in subsequent turns.

Increasing batch size B from 1 to 5 substantially improves overall performance, while further increasing it to 8 results in a performance drop. We conjecture that this decline may be due to the LLM failing to correctly identify the best and worst solutions when the number of candidate solutions grows, highlighting an open problem of how to perform this selection optimally.

Fig. 3 compares 3D Scaling with baseline methods across different parameter settings. The results show that 3D Scaling effectively leverages the reasoning capabilities of the LLM, achieving a maximum average accuracy of 73.3%. We also report the outcome of applying 3D Scaling with human judgment under the setting $C = 32768, B = 5, T = 3$, where the score rises to 86.7%, demonstrating the potential benefits of incorporating human preferences.

Insights: Performance improves along all three scaling dimensions up to a point, but context scaling plateaus, batch scaling has an optimal value, and turn scaling also saturates; whether additional dimensions could further enhance reasoning remains an open question.

4.3 RESULTS ON BENCHMARK TASKS

In this subsection, we present 3D scaling experiments with selection feedback from both LLM and human judges on three challenging benchmarks, using a setting of $B = 5, T = 3$. For statistical reliability, we conducted 10 independent trials for the single-run CoT baseline and 5 trials for each of the other comparative methods. To ensure fairness, batch scaling was configured to generate 15 solutions per trial, thereby matching the total token budget of the 3D scaling setup. We report the average of the final scores across different trials and the standard deviation.

4.3.1 MATH OLYMPICS

The performance of different TTS methods on all six problems in IMO 2025 (International Mathematical Olympiad (2025)) is summarized in Table 1. The experimental results reveal several key ob-

324 observations. The single-response **Context Scaling** approach achieved moderate performance. Anal-
 325 ysis of the responses indicates that while the model can produce reasonable answers over multiple
 326 runs, it often generates incomplete or partially valid reasoning. **Batch Scaling (Vote)** and **Turn**
 327 **Scaling (Reflection)** improve accuracy over the context scaling baseline by scaling along individ-
 328 ual dimensions. However, both methods reach saturation when the scale increases to 15, and the
 329 model’s ability to produce fully complete reasoning remains limited. The fully automated iterative
 330 refinement approach, **3D Scaling (LLM Judge)**, demonstrates competitive performance, achieving
 331 higher accuracy than the baseline scaling methods. This suggests that scaling across multiple dimen-
 332 sions can overcome the limitations of single-dimension scaling. Furthermore, applying **3D Scaling**
 333 **(Human Judge)** leads to substantial improvements over all baselines, achieving the best overall
 334 performance. Incorporating human feedback addresses the LLM’s tendency to produce incomplete
 335 reasoning, enabling it to generate solutions with fully correct reasoning through iterative refinement.

336 Table 1: Average accuracy of different test-time scaling methods on IMO 2025. The values in
 337 parentheses represent the standard deviation.

Method	IMO1	IMO2	IMO3	IMO4	IMO5	IMO6
Context Scaling	10%(30%)	20%(40%)	30%(46%)	80%(40%)	40%(49%)	0%(0%)
Batch Scaling (Vote)	20%(40%)	/	40%(49%)	100% (0%)	20%(40%)	0%(0%)
Turn Scaling (Reflection)	60%(49%)	40%(49%)	60%(49%)	100% (0%)	60%(49%)	0%(0%)
3D Scaling (LLM Judge)	60%(49%)	40%(49%)	80%(40%)	100% (0%)	80% (40%)	0%(0%)
3D Scaling (Human Judge)	100% (0%)	60% (49%)	100% (0%)	100% (0%)	60%(49%)	0%(0%)

345 4.3.2 PHYSICS OLYMPICS

346 The performance of different TTS methods on all six problems in CPHO 2022 is summarized in
 347 Table 2. The results on physics competition problems demonstrate a trend consistent with that ob-
 348 served in mathematical competitions: **3D Scaling (Human Judge)** achieves the highest accuracy,
 349 followed by **3D Scaling (LLM Judge)**, **Batch Scaling (Vote)**, and finally the single-response **Con-**
 350 **text Scaling**. Extra analysis is provided in Appendix A.3.

351 Table 2: Average accuracy of different reasoning methods on CPHO 2022. The values in parentheses
 352 represent the standard deviation.

Method	CPHO1	CPHO2	CPHO3	CPHO4	CPHO5	CPHO6
Context Scaling	70%(46%)	100% (0%)	0%(0%)	0%(0%)	40%(49%)	100% (0%)
Batch Scaling (Vote)	80%(40%)	100% (0%)	0%(0%)	0%(0%)	80%(40%)	100% (0%)
3D Scaling (LLM Judge)	100% (0%)	100% (0%)	20% (40%)	0%(0%)	60%(49%)	100% (0%)
3D Scaling (Human Judge)	100% (0%)	100% (0%)	20% (40%)	0%(0%)	100% (0%)	100% (0%)

360 4.3.3 CODING

361 The test results of different TTS methods on IOI 2025 are presented in Table 3. Among the six
 362 problems in IOI 2025, the fifth problem is a communication task; since the backbone model cannot
 363 access the submission API, its performance on this task is unsatisfactory, and we therefore exclude
 364 it from evaluation. The second problem consists of two parts worth 70 and 30 points, respectively,
 365 while all remaining problems are scored out of 100 points.

366 Table 3: Average scores of different test-time scaling methods on IOI 2025. The values in parenthe-
 367 ses represent the standard deviation.

Method	IOI1	IOI2-Part1	IOI2-Part2	IOI3	IOI4	IOI6	Average Score
Context Scaling	2.3(2.68)	24.7(23.7)	2.79(1.80)	8.4(4.15)	25.6(19.3)	13.2(17.3)	12.82
Batch Scaling (Best of N)	9.8(5.6)	53.4(16.8)	5.74(0.90)	32(21.42)	57.2(12.4)	26.6(5.2)	30.79
3D Scaling (LLM Judge)	3.4(2.24)	38(18.8)	5.13(1.43)	10.2(5.42)	32(23.6)	21.2(5.63)	18.65
3D Scaling (Human Judge)	12.6(6.85)	70(0)	5.13(1.15)	25.6(16.4)	65.4(1.34)	42.8(14.8)	36.59

370 The results indicate that 3D Scaling can substantially enhance coding performance through human
 371 feedback. On IOI problems, LLMs often struggle to generate fully correct solutions in a zero-
 372 shot setting. Consequently, **Context Scaling** typically solves only a subset of tasks and sometimes
 373 contains errors in complexity analysis. Because the algorithms initially generated by the LLM vary
 374

378 significantly across IOI problems, we adopt **Batch Scaling (Best of N)** to demonstrate the upper
379 bound of Batch Scaling.

380
381 When feedback from a human is incorporated, 3D Scaling demonstrates the best performance. Even
382 when all early solutions are incorrect, the **3D Scaling (LLM Judge)** can identify issues in the
383 code and iteratively refine them, enabling it to solve more tasks and achieve higher scores. Across
384 nearly all experiments, this approach produces final scores that surpass the best first-round solutions,
385 achieving an average improvement of approximately 18.8% over the Batch Scaling baseline.

386 While the **3D Scaling (LLM Judge)** is less precise than a **3D Scaling (Human Judge)** on chal-
387 lenging problems, it consistently outperforms the **Context Scaling** approach and achieves compa-
388 rable overall performance. These results highlight the feasibility and effectiveness of auto-feedback
389 mechanisms for improving code generation, even in the absence of human feedback.

390 4.4 EXPERIMENTS ON INNOVATIVE TASKS

391
392 In this section, we evaluate the effects of human feedback on several robotics reinforcement learning
393 tasks using the GPU-based IsaacGym framework (Makoviychuk et al. (2021)), including *Cartpole*,
394 *BallBalance*, *Quadcopter*, *Ant*, *Humanoid*, *ShadowHand*, and *AllegroHand*, along with a challeng-
395 ing and innovative new task, *HumanoidJump*, defined as “making a humanoid jump like a real
396 human.”, which is an open-ended challenge without gold-standard answers.

397 We employed GPT-4o as the backbone model. The model was prompted to generate task-specific
398 reward functions, which were then used to train agents in the simulator. In these tasks, we employed
399 settings of $B = 6$ and $T = 5$ for **3D Scaling (Human Judge)**. In each iteration, the evaluators
400 selected the best and worst reward functions based on behavior videos of the agents trained with
401 these reward functions. Details about this process are provided in Appendix A.4. We also report the
402 results with the **Context Scaling** and **Batch Scaling (Best of N)**.

403 In each turn, in addition to providing preference feedback, we also generate automatic feedback
404 with LLM, which is combined with human preferences as the feedback prompt for the next round to
405 assist the LLM in refinement. The automatic feedback consists of the following three components:

- 407 • **Evaluation of reward functions:** The component values that make up the good and bad
408 reward functions are obtained from the environment during training and provided to the
409 LLM. This helps the LLM assess the usefulness of different parts of the reward function by
410 comparing the two.
- 411 • **Differences between historical reward functions:** We employed GPT-4o to analyze the
412 differences between the historically best reward functions from each iteration. These dif-
413 ferences were then provided to the generator LLM to assist in refining the reward function.
- 414 • **Reward trace of historical reward functions:** The reward trace, consisting of the values
415 of the good reward functions during training from all prior iterations, is provided to the
416 LLM. This reward trace enables the LLM to evaluate how well the agent is actually able to
417 optimize those reward components.

418 4.4.1 TASK METRIC

419
420 For evaluation, we used the reward function in a PPO (Schulman et al., 2017) training loop following
421 the original setting in IsaacGym, and reported the average task score, measured by the expert-written
422 task metrics across multiple experiments, as the ground truth rewards for each method. The details
423 of the task metrics are provided in Appendix B.2. For the *HumanoidJump* task, since designing a
424 reward metric is challenging, we adopt human votes for quantitative evaluation instead, which is
425 detailed in Sec. 4.4.3

426 4.4.2 ISAACGYM TASKS RESULTS

427
428 For each environment, we conducted five runs per method and reported the average ground-truth
429 rewards in Table 4, while ensuring that **Batch Scaling (Best of N)** and **3D Scaling (Human Judge)**
430 used the same total token budget. As observed, **3D Scaling (Human Judge)** significantly outper-
431 forms **Batch Scaling (Best of N)** in 3 out of 5 tasks, achieving an average improvement of 18.4%.

Table 4: Average ground truth rewards of different test-time scaling methods on IsaacGym Tasks. The values in parentheses represent the standard deviation.

	Cart.	Ball.	Quad.	Ant	Human.	Shadow	Allegro
Context Scaling	499(0)	499(0)	-0.356(0.29)	5.262(2.49)	6.157(0.86)	6.605(2.95)	15.500(9.34)
Batch Scaling (Best of N)	499(0)	499(0)	-0.0410(0.32)	9.350(2.34)	8.306(1.63)	9.476(2.44)	23.876(7.91)
3D Scaling (Human Judge)	499(0)	499(0)	-0.0183(0.29)	11.142(0.37)	8.392(0.53)	10.740(0.92)	24.134(6.52)

In addition, we conducted another set of experiments with a proxy judge and analyzed performance improvements across turns, as detailed in Appendix B.

4.4.3 HUMANOIDJUMP TASK RESULTS

Without human feedback, the most common behavior observed in this task was what we refer to as a “leg-lift jump.” In this behavior, the robot lifts only one leg to shift its center of mass and attempts a jump, while the opposite leg pushes off the ground to achieve lift. Although this behavior satisfies the minimal definition of a jump—reaching a certain distance above the ground, it does not align with our expectation of a human-like jump.

In contrast, when real human preferences were incorporated, the outcomes were notably different. The volunteers judged the overall quality of the humanoid’s jump behavior rather than relying solely on the metric of leaving the ground. The results show that human feedback can effectively guide the humanoid toward more natural, human-like jumps by favoring behaviors that, although not initially optimal, exhibited promising movement patterns. After six iterations, the humanoid displayed more sophisticated behaviors, such as bending both legs and lowering the upper body to shift its center of mass—motions much closer to a real human jump. A detailed analysis and illustrative images are provided in Appendix C.

For quantitative evaluation, we adopt human votes for the quantitative evaluation on HumanoidJump task. As a baseline, we use the **Batch Scaling (Best of N)**, which generates 30 reward functions with the same total token budget as **3D Scaling (Human Judge)**. 20 volunteers were recruited to compare the performance of the two methods. Each volunteer indicated their preference between two videos presented in random order—one generated by **3D Scaling (Human Judge)** and the other by **Batch Scaling (Best of N)**. As shown in Table 5, 17 out of 20 participants preferred the **3D Scaling (Human Judge)** agent, demonstrating that **3D Scaling (Human Judge)** produces behaviors more aligned with human preferences.

Table 5: Human Preferences over different agents.

Method	Vote
Batch Scaling (Best of N)	3/20
3D Scaling (Human Judge)	17/20

5 CONCLUSION

In this work, we revisited test-time enhancement techniques for reasoning models from the perspective of scaling laws. By unifying existing approaches under the framework of multi-dimensional test-time scaling, we identified three orthogonal axes—context, batch, and turn—each of which independently exhibits a clear scaling law. Building on this observation, we introduced 3D test-time scaling, which integrates all three dimensions to substantially extend the effective capacity of test-time compute. Our experiments demonstrated that this unified framework not only improves reasoning accuracy on challenging benchmarks such as IOI, IMO, and CPHO, but also naturally supports a human-in-the-loop paradigm that further amplifies model performance. Moreover, we showed that the same principles can be applied to embodied learning, enabling reasoning models to discover novel behaviors for humanoid robot control.

Despite these advances, important open questions remain. While our study has revealed three fundamental scaling dimensions, the capacity of test-time compute is still bottlenecked by architectural and computational constraints. It remains unclear whether additional dimensions of scaling exist that could further unlock the reasoning potential of large models. Exploring such new axes—beyond context, batch, and turn—represents an exciting direction for future research.

REFERENCES

- 486
487
488 Pranjali Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with
489 reinforcement learning. In *Second Conference on Language Modeling, 2025*. URL <https://openreview.net/forum?id=4jdIxXBNve>.
490
- 491 Chinese Physics Olympiad. <http://cpho.pku.edu.cn/>, 2022.
492
- 493 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
494 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
495 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 496 Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
497 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, et al. Gemini 2.5: Pushing the frontier with
498 advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025.
499 URL <https://arxiv.org/abs/2507.06261>.
- 500 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,
501 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms
502 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
503
- 504 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
505 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan,
506 Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon
507 Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training
508 compute-optimal large language models, 2022. URL <https://arxiv.org/abs/2203.15556>.
- 509 Yichen Huang and Lin F. Yang. Gemini 2.5 pro capable of winning gold at imo 2025, 2025. URL
510 <https://arxiv.org/abs/2507.15855>.
- 511 International Mathematical Olympiad. <https://www.imo-official.org/>, 2025.
512
- 513 International Olympiad in Informatics. <https://ioinformatics.org/>, 2025.
514
- 515 International Physics Olympiad. <https://www.ipho-new.org/>.
- 516 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child,
517 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language
518 models, 2020. URL <https://arxiv.org/abs/2001.08361>.
- 519 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri
520 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement
521 with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023.
522
- 523 Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin,
524 David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High
525 performance gpu-based physics simulation for robot learning, 2021. URL <https://arxiv.org/abs/2108.10470>.
526
- 527 Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke
528 Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time
529 scaling, 2025. URL <https://arxiv.org/abs/2501.19393>.
- 530 Xuefei Ning, Zinan Lin, Zixuan Zhou, Zifu Wang, Huazhong Yang, and Yu Wang. Skeleton-of-
531 thought: Prompting llms for efficient parallel generation. *arXiv preprint arXiv:2307.15337*, 2023.
532
- 533 OpenAI. Learning to reason with llms. <https://openai.com/index/learning-to-reason-with-llms>, 2024.
534
- 535 Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song,
536 John Aslanides, et al. Scaling language models: Methods, analysis insights from training gopher,
537 2022. URL <https://arxiv.org/abs/2112.11446>.
- 538 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
539 optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.

540 Jingzhe Shi, Qinwei Ma, Hongyi Liu, Hang Zhao, Jeng-Neng Hwang, and Lei Li. Explaining
541 context length scaling and bounds for language models, 2025. URL [https://arxiv.org/abs/
542 2502.01481](https://arxiv.org/abs/2502.01481).

543 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
544 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing
545 Systems*, 36:8634–8652, 2023.

547 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
548 can be more effective than scaling model parameters, 2024. URL [https://arxiv.org/abs/
549 2408.03314](https://arxiv.org/abs/2408.03314).

550 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
551 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models,
552 2023a. URL <https://arxiv.org/abs/2203.11171>.

553 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
554 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models,
555 2023b. URL <https://arxiv.org/abs/2203.11171>.

557 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc
558 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models,
559 2023. URL <https://arxiv.org/abs/2201.11903>.

560 Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck, and Yiming Yang. Inference scaling laws:
561 An empirical analysis of compute-optimal inference for problem-solving with language models,
562 2025. URL <https://arxiv.org/abs/2408.00724>.

563 Yuxi Xie, Anirudh Goyal, Wenyue Zheng, Min-Yen Kan, Timothy P Lillicrap, Kenji Kawaguchi,
564 and Michael Shieh. Monte carlo tree search boosts reasoning via iterative preference learning.
565 *arXiv preprint arXiv:2405.00451*, 2024.

566 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik
567 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Ad-
568 vances in neural information processing systems*, 36:11809–11822, 2023.

569 Dan Zhang, Sining Zhou, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. Rest-mcts*: Llm
570 self-training via process reward guided tree search. *Advances in Neural Information Processing
571 Systems*, 37:64735–64772, 2024.

572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

594 A EXPERIMENTS DETAILS

595 A.1 PREFERENCE FEEDBACK

596 We evaluate 3D Scaling under two feedback mechanisms:

- 597 • **LLM Judge:** In each iteration, the model generates B candidate responses, and Gemini-2.5-Pro selects the best and worst solutions through pairwise comparisons. .
- 600 • **Human Judge:** Three volunteers, each of whom has won a gold medal in a national Olympiad in mathematics, physics, or informatics, serve as evaluators. In each iteration, they identify the best and worst responses among the candidates. They do not have access to the official reference solutions and provide no additional explanations, with only the selected results fed back to the model. For the coding benchmark, the evaluators additionally compile and run the generated code, testing it against problem-specific subtasks and constraints.

609 A.2 BENCHMARK CHOICE

610 The CPHO dataset was selected over the IPHO dataset(International Physics Olympiad primarily because IPHO problems are typically decomposed into a large number of weakly related sub-questions (e.g., 20 per problem), making it computationally expensive to evaluate the quality of each individual response. In contrast, CPHO problems contain fewer sub-questions (e.g., 5 per problem) and exhibit strong logical coherence across all parts of a given problem. As a result, the correctness of the last sub-question can serve as a reliable indicator of whether the model has successfully solved the entire problem.

618 A.3 EXTRA ANALYSIS ON CPHO

619 Regarding specific accuracy rates, the LLM exhibits the following characteristics when solving physics competition problems:

- 620 1. **Difference in Evaluation between Physics and Mathematics Problems:** Unlike mathematics problems, where the final answer may be relatively straightforward to conjecture while the reasoning process can be highly complex, physics problems typically feature a final answer that is difficult to obtain. Consequently, if a correct final answer is produced, it generally indicates a valid reasoning process. As a result, the evaluation of physics solutions relies almost entirely on the correctness of the final answer.
- 628 2. **Instability of Final Answers during Self-Improvement:** In contrast to mathematics problems, during self-improvement iterations, the model exhibits a higher tendency to alter the final answer, reflecting greater uncertainty or refinement in the solution process for physics questions.

633 A.4 HUMAN EVALUATION DETAILS

634 We conducted human-in-the-loop experiments with human participants. During each iteration, human evaluators select the optimal and most deficient solutions among these candidates based on whether they satisfy the task objectives and whether they can be further improved.

635 The human evaluators are three volunteers, each of whom has won a gold medal in a national-level Olympics competition in mathematics, physics, or informatics. Only the best and worst solutions themselves are fed back to the LLM to guide further self-refinement; evaluators do not provide any information about the reasons for their choices or about bugs in these responses.

636 During the human evaluation process, the annotators were provided with the standard answers to the mathematics and physics problems. The evaluation protocol was as follows: annotators first assessed whether the final answer provided in the model’s response was correct. Only if the final answer was correct did they proceed to evaluate the reasonableness of the key steps within the solution process.

641 Given the strong interdependence between subproblems within the CPHO physics problems, we manually identified and tagged the final logical step of each problem as a *key subproblem*. In the

system prompt, the model was explicitly instructed to present its response to this key subproblem at the very beginning of its overall reply. This design allows human annotators to quickly gauge the problem’s overall correctness; if the answer to the key subproblem is correct, it serves as a strong indicator that the entire problem has likely been solved correctly.

For IOI tasks, the evaluators additionally compile and run the code generated by the model, testing it against test cases that satisfy problem-specific subtasks and constraints.

B EXTRA EXPERIMENTS ON ISAACGYM TASKS

In this section, we discussed the details of experiments on IsaacGym tasks.

B.1 ENVIRONMENT DETAILS

In Table 6, we present the observation and action dimensions, along with the task description and task metrics for 9 tasks in IsaacGym.

Environment (obs dim, action dim)
Task Description
Task Metric
Cartpole (4, 1)
To balance a pole on a cart so that the pole stays upright
<i>duration</i>
Quadcopter (21, 12)
To make the quadcopter reach and hover near a fixed position
<i>-cur_dist</i>
FrankaCabinet (23, 9)
To open the cabinet door
<i>1 if cabinet_pos > 0.39</i>
Anymal (48, 12)
To make the quadruped follow randomly chosen x, y, and yaw target velocities
<i>-(linvel_error + angvel_error)</i>
BallBalance (48, 12)
To keep the ball on the table top without falling
<i>duration</i>
Ant (60, 8)
To make the ant run forward as fast as possible
<i>cur_dist - prev_dist</i>
AllegroHand (88, 16)
To make the hand spin the object to a target orientation
<i>number of consecutive successes where current success is 1 if rot_dist < 0.1</i>
Humanoid (108, 21)
To make the humanoid run as fast as possible
<i>cur_dist - prev_dist</i>
ShadowHand (211, 20)
To make the shadow hand spin the object to a target orientation
<i>number of consecutive successes where current success is 1 if rot_dist < 0.1</i>

Table 6: Details of IsaacGym Tasks.

B.2 TASK METRICS

we employed the average of the sparse rewards across parallel environments as the task metrics, following the original setting in IsaacGym.

Table 7: Average FTS of different test-time scaling methods on IsaacGym Tasks. The values in parentheses represent the standard deviation.

	Cart.	Ball.	Quad.	Ant	Human.	Shadow	Allegro
Context-Scaling	499(0)	499(0)	-0.356(0.29)	5.262(2.49)	6.157(0.86)	6.605(2.95)	15.500(9.34)
Batch-Scaling(Best of all)	499(0)	499(0)	-0.0410(0.32)	9.350(2.34)	8.306(1.63)	9.476(2.44)	23.876(7.91)
3D Scaling(Proxy Judge)	499(0)	499(0)	-0.0195(0.09)	12.04(1.69)	9.227(0.93)	13.231(1.88)	25.030(3.721)
3D Scaling(Human Judge)	499(0)	499(0)	-0.0183(0.29)	11.142(0.37)	8.392(0.53)	10.740(0.92)	24.134(6.52)

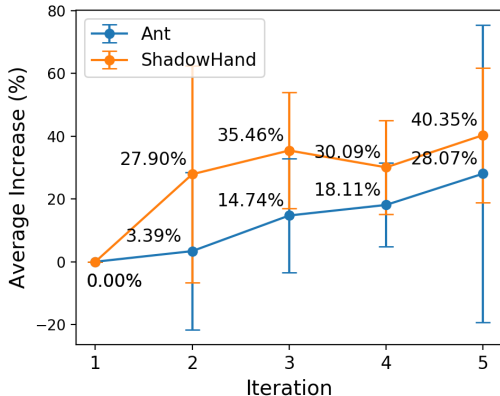


Figure 4: Average improvement of the Reward Task Score (RTS) compared with the first iteration in 3D scaling-Proxy Judge for the Ant and ShadowHand tasks, demonstrating the method’s effectiveness in refining reward functions.

To assess the generated reward function for each RL run, we take the maximum task metric value sampled at fixed intervals, marked as *task score of reward function* (RTS). In each iteration, 3D scaling generates $B = 6$ RL runs and selects the highest RTS as the result for that iteration. 3D scaling performs $T = 5$ iterations and then selects the highest RTS from these iterations as the *task score* (TS) for each experiment. Due to the inherent randomness of LLMs, we run 5 experiments for all methods, and report the highest TS as the *final task score* (FTS) for each approach. A higher FTS indicates better performance across all tasks.

B.3 3D SCALING WITH PROXY JUDGE

In IsaacGym tasks, it is difficult for an LLM to evaluate the quality of reward functions from videos as humans do. To address this, we use human-designed expert rewards as a proxy for human preference, enabling rapid and quantitative evaluation of our approach. This proxy represents a noise-free case that is likely easier than real human trials. Importantly, these human-designed rewards are used solely to automate sample selection and are never included in the prompts sent to the LLM; the LLM never observes the functional form of the ground-truth rewards nor receives any values from them. The results are referred to as ‘3D scaling(Proxy Judge)’ in the tables.

B.4 IMPROVEMENT ANALYSIS

We provided the final average FTS with an extra variant in Table 7. As observed, on average, “3D Scaling(Proxy Judge)” achieves a 27.4% improvement over Batch-Scaling. We can also observe that 3D Scaling exhibits lower variance than Batch Scaling, indicating more stable reward learning behavior.

While it is possible that the LLMs could generate an optimal reward function in a zero-shot manner, the primary focus of our analysis is not solely on absolute performance values. Rather, we emphasize whether 3D Scaling is capable of enhancing performance through the iterative incorporation of preferences. We calculated the average RTS improvement compared to the first iteration for the two tasks with the largest improvements compared with Batch-Scaling(Best of N), *Ant*, and *Shadow-Hand*. As shown in Fig. 4, RTS demonstrates improved performance after multiple iterations (e.g., 5 vs. 1), highlighting its effectiveness in refining reward functions.

B.5 PSEUDOCODE

The full pseudocode of 3D Scaling on embodied AI tasks is listed in Algo. 1.

Algorithm 1: 3D Scaling

```

756 Input: # iterations  $N$ , # samples in each iterations  $K$ , environment Env, coding LLM  $LLM_{RF}$ ,
757         difference LLM  $LLM_{Diff}$ 
758
759
760 Function Feedback(Env, RF):
761     return The values of each component that make up RF during the training process in Env
762
763 Function History(RFlist, Env,  $LLM_{Diff}$ ):
764     HistoryFeedback  $\leftarrow$  ""
765     for  $i \leftarrow 1$  to  $len(RFlist) - 1$  do
766         // The reward trace of historical reward functions
767         HistoryFeedback  $\leftarrow$  HistoryFeedback + Feedback(Env, RFlist[ $i - 1$ ])
768         // The differences between historical reward functions
769         HistoryFeedback  $\leftarrow$ 
770             HistoryFeedback +  $LLM_{Diff}$ (DifferencePrompt + RFlist[ $i$ ] + RFlist[ $i - 1$ ])
771     end
772     return HistoryFeedback
773
774 // Initialize the prompt containing the environment context and task description
775 Prompt  $\leftarrow$  InitializePrompt
776 RFlist  $\leftarrow$  []
777 for  $i \leftarrow 1$  to  $N$  do
778     RF $_1, \dots, RF_K \leftarrow LLM_{RF}$ (Prompt,  $K$ )
779     while any of RF $_1, \dots, RF_K$  is not executable do
780          $j_1, \dots, j_{K'} \leftarrow$  Index of non-executable reward functions
781         // Regenerate non-executable reward functions
782         RF $_{j_1}, \dots, RF_{j_{K'}} \leftarrow LLM_{RF}$ (Prompt,  $K'$ )
783     end
784     // Render videos for sampled reward functions
785     Video $_1, \dots, Video_K \leftarrow$  Render(Env, RF $_1$ ),  $\dots$ , Render(Env, RF $_K$ )
786     // Human selects the most preferred and least preferred videos
787      $G, B \leftarrow$  Human(Video $_1, \dots, Video_K$ )
788     GoodRF, BadRF  $\leftarrow$  RF $_G, RF_B$ 
789     RFlist.append(GoodRF)
790     // Update prompt for feedback
791     Prompt  $\leftarrow$ 
792         GoodRF + Feedback(Env, GoodRF) + BadRF + Feedback(Env, BadRF) + PreferencePrompt
793     Prompt  $\leftarrow$  Prompt + History(RFlist, Env,  $LLM_{Diff}$ )
794 end

```

B.6 EXAMPLE

We use a trial of the *Humanoid* task to illustrate how 3D Scaling progressively generated improved reward functions over successive iterations. The task description is “to make the humanoid run as fast as possible”. Throughout five iterations, adjustments were made to the penalty terms and reward weightings. In the first iteration, the total reward was calculated as $0.5 \times \text{speed_reward} + 0.25 \times \text{deviation_reward} + 0.25 \times \text{action_reward}$, yielding an RTS of 5.803. The speed reward and deviation reward motivate the humanoid to run fast, while the action reward promotes smoother motion. In the second iteration, the weight of the speed reward was increased to 0.6, while the weights for deviation and action rewards were adjusted to 0.2 each, improving the RTS to 6.113. In the third iteration, the action penalty was raised and the reward weights were further modified to $0.7 \times \text{speed_reward}$, $0.15 \times \text{deviation_reward}$, and $0.15 \times \text{action_reward}$, resulting in an RTS of 7.915. During the fourth iteration, the deviation penalty was reduced to 0.35 and the action penalty was lowered, with the reward weights set to 0.8, 0.1, and 0.1 for speed, deviation, and action rewards, respectively. This change led to an RTS of 8.125. Finally, in the fifth iteration, an additional upright reward term was incorporated, with the total reward calculated as $0.7 \times \text{speed_reward} + 0.1 \times \text{deviation_reward} + 0.1 \times$

810 action_reward+0.1×upright_reward. This adjustment produced the highest RTS of 8.232, allowing
 811 3D Scaling to generate reward functions that were more effectively aligned with the task description.
 812 Below are the specific reward functions produced at each iteration during one experiment.
 813

814 Humanoid Task: Reward Function with highest RTS (5.803) of Iteration 1

```

815
816 1 def compute_reward(root_states: torch.Tensor, actions: torch.Tensor
817 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
818     2 velocity = root_states[:, 7:10]
819     3 forward_velocity = velocity[:, 0]
820     4 target_velocity = 5.0
821     5 deviation_penalty = 0.5
822     6 action_penalty = 0.1
823
824     7
825     8 # Measure how fast the humanoid is going
826     9 speed_reward = torch.exp((forward_velocity - target_velocity))
827
828     0
829     1 # Penalize deviation from the target velocity
830     2 deviation = torch.abs(forward_velocity - target_velocity)
831     3 deviation_reward = torch.exp(-deviation_penalty * deviation)
832
833     4
834     5 # Penalize actions, to induce smoother motion
835     6 sum_actions = torch.sum(actions ** 2, dim=1)
836     7 action_reward = torch.exp(-action_penalty * sum_actions)
837
838     8
839     9 rewards = {'speed_reward': speed_reward, 'deviation_reward':
840               deviation_reward, 'action_reward': action_reward }
841
842     0
843     1 total_reward = 0.5 * speed_reward + 0.25 * deviation_reward +
844               0.25 * action_reward
845     2 total_reward = total_reward.to(deviation.device)
846
847     3
848     4 return total_reward, rewards
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

```

841 Humanoid Task: Reward Function with highest RTS (6.113) of Iteration 2

```

842
843 1 def compute_reward(root_states: torch.Tensor, actions: torch.Tensor
844 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
845     2 velocity = root_states[:, 7:10]
846     3 forward_velocity = velocity[:, 0]
847     4 target_velocity = 5.0
848     5 deviation_penalty = 0.5
849     6 action_penalty = 0.1
850
851     7
852     8 # Measure how fast the humanoid is going
853     9 speed_reward = torch.exp((forward_velocity - target_velocity))
854
855     0 # Penalize deviation from the target velocity
856     1 deviation = torch.abs(forward_velocity - target_velocity)
857     2 deviation_reward = torch.exp(-deviation_penalty * deviation)
858
859     3 # Penalize actions, to induce smoother motion
860     4 sum_actions = torch.sum(actions ** 2, dim=1)
861     5 action_reward = torch.exp(-action_penalty * sum_actions)
862
863     6
864     7 rewards = {'speed_reward': speed_reward, 'deviation_reward':
865               deviation_reward, 'action_reward': action_reward }
866
867     8
868     9 # Adjusted weights in total reward calculation
869     0 total_reward = 0.6 * speed_reward + 0.2 * deviation_reward +
870               0.2 * action_reward # Changed weights
871     1 total_reward = total_reward.to(deviation.device)
872
873
874
875
876
877
878
879
880
881
882
883

```

```

864
865
866     return total_reward, rewards
867

```

Humanoid Task: Reward Function with highest RTS (7.915) of Iteration 3

```

870
871 1 def compute_reward(root_states: torch.Tensor, actions: torch.Tensor
872 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
873     velocity = root_states[:, 7:10]
874     forward_velocity = velocity[:, 0]
875     target_velocity = 5.0
876     deviation_penalty = 0.5
877     action_penalty = 0.15 # Increased action penalty
878
879     # Measure how fast the humanoid is going
880     speed_reward = torch.exp((forward_velocity - target_velocity))
881     # Penalize deviation from the target velocity
882     deviation = torch.abs(forward_velocity - target_velocity)
883     deviation_reward = torch.exp(-deviation_penalty * deviation)
884     # Penalize actions, to induce smoother motion
885     sum_actions = torch.sum(actions ** 2, dim=1)
886     action_reward = torch.exp(-action_penalty * sum_actions)
887
888     rewards = {'speed_reward': speed_reward, 'deviation_reward':
889               deviation_reward, 'action_reward': action_reward }
890
891     # Adjusted weights in total reward calculation
892     total_reward = 0.7 * speed_reward + 0.15 * deviation_reward +
893                 0.15 * action_reward # Changed weights
894     total_reward = total_reward.to(deviation.device)
895
896     return total_reward, rewards
897

```

Humanoid Task: Reward Function with highest RTS (8.125) of Iteration 4

```

898 1 def compute_reward(root_states: torch.Tensor, actions: torch.Tensor
899 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
900     velocity = root_states[:, 7:10]
901     forward_velocity = velocity[:, 0]
902     target_velocity = 5.0
903     deviation_penalty = 0.35 # decrease penalty for deviation
904     action_penalty = 0.1 # decrease penalty for action
905
906     # Measure how fast the humanoid is going
907     speed_reward = torch.exp((forward_velocity - target_velocity))
908     # Penalize deviation from the target velocity
909     deviation = torch.abs(forward_velocity - target_velocity)
910     deviation_reward = torch.exp(-deviation_penalty * deviation)
911     # Penalize actions, to induce smoother motion
912     sum_actions = torch.sum(actions ** 2, dim=1)
913     action_reward = torch.exp(-action_penalty * sum_actions)
914
915     rewards = {'speed_reward': speed_reward, 'deviation_reward':
916               deviation_reward, 'action_reward': action_reward }
917
918     # Adjust the weights in total reward calculation
919     total_reward = 0.8 * speed_reward + 0.1 * deviation_reward +
920                 0.1 * action_reward # Changed weights
921     total_reward = total_reward.to(deviation.device)
922

```

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954

```

2
3     return total_reward, rewards

```

```

Humanoid Task: Reward Function with highest RTS (8.233) of Iteration 5

1  def compute_reward(root_states: torch.Tensor, actions: torch.Tensor
2     , up_vec: torch.Tensor) -> Tuple[torch.Tensor, Dict[str, torch.
3     Tensor]]:
4     velocity = root_states[:, 7:10]
5     forward_velocity = velocity[:, 0]
6     target_velocity = 5.0
7     deviation_penalty = 0.35 # Reduced deviation penalty
8     action_penalty = 0.1 # Decreased action penalty
9
10    # Measure how fast the humanoid is going
11    speed_reward = torch.exp((forward_velocity - target_velocity))
12    # Penalize deviation from the target velocity
13    deviation = torch.abs(forward_velocity - target_velocity)
14    deviation_reward = torch.exp(-deviation_penalty * deviation)
15    # Penalize actions, to induce smoother motion
16    sum_actions = torch.sum(actions ** 2, dim=1)
17    action_reward = torch.exp(-action_penalty * sum_actions)
18    # Reward for maintaining an upright position
19    upright_penalty = 1.0 # New upright penalty for the humanoid
20    upright_reward = torch.exp(-upright_penalty * (1 - up_vec[:,
21        2])) # Added upright reward
22
23    rewards = {'speed_reward': speed_reward, 'deviation_reward':
24        deviation_reward, 'action_reward': action_reward, '
25        upright_reward': upright_reward }
26
27    # Adjusted weights in total reward calculation
28    total_reward = 0.7 * speed_reward + 0.1 * deviation_reward +
29        0.1 * action_reward + 0.1 * upright_reward # Added upright
30        reward to total
31    total_reward = total_reward.to(deviation.device)
32
33    return total_reward, rewards

```

C EXPERIMENTS ON HUMANOIDJUMP TASK

In this section, we analyzed the effectiveness of human feedback in HumanoidJump task through analysis on the behaviors of the agents.

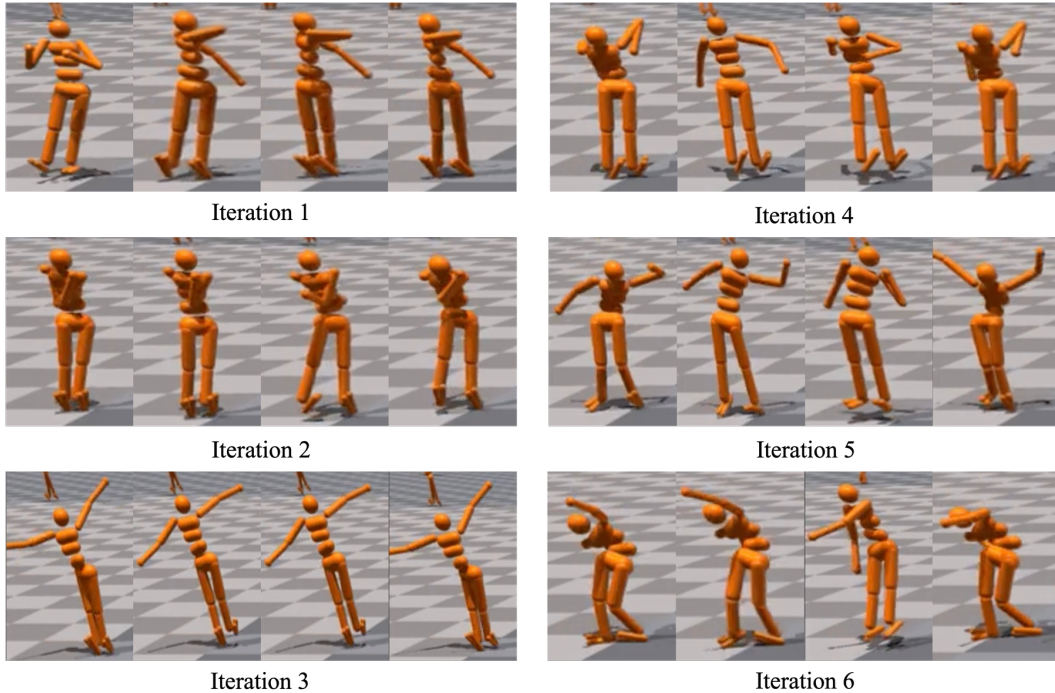
960
961
962
963
964
965
966
967
968
969
970
971



Figure 5: A common behavior.

The most common behavior observed in this task, as illustrated in Fig. 5, is what we refer to as the “leg-lift jump.” This behavior involves initially lifting one leg to raise the center of mass, followed

972 by the opposite leg pushing off the ground to achieve lift. The previously lifted leg is then lowered
 973 to extend airtime. Various adjustments of the center of mass with the lifted leg were also noted.
 974 This behavior meets the minimal metric of a jump: achieving a certain distance off the ground. If
 975 feedback were provided based solely on this minimal metric, the “leg-lift jump” would likely be
 976 selected as a candidate reward function. However, Such candidates show limited improvement in
 977 subsequent iterations, failing to evolve into more human-like jumping behaviors.



1001 Figure 6: The humanoid learns a human-like jump by bending legs and lowering the upper body to
 1002 shift the center of mass in a trial of human-in-the-loop experiments. Note that both legs are used to
 1003 jump, and the agent bends at the hips.

1005 Conversely, when real human preferences were used to guide the task, the results were notably dif-
 1006 ferent. The volunteer judged the overall quality of the humanoid’s jump behavior instead of just
 1007 the metric of leaving the ground. Fig. 6 illustrates that the volunteer successfully guided the hu-
 1008 manoid towards a more human-like jump by selecting behaviors that, while initially not optimal,
 1009 displayed promising movement patterns. In the first iteration, “leg-lift jump” was not selected de-
 1010 spite the humanoid jumping off the ground. Instead, a video where the humanoid appears to attempt
 1011 a jump using both legs, without leaving the ground, was chosen. By the fifth and sixth iterations,
 1012 the humanoid demonstrated more sophisticated behaviors, such as bending both legs and lowering
 1013 the upper body to shift the center of mass, behaviors that are much more akin to a real human jump.
 1014 The videos can be found on our website.

1015 D FULL PROMPTS

1016 D.1 FULL PROMPTS ON EMBODIED AI TASKS

1017 The prompts used in 3D Scaling for synthesizing reward functions in Embodied AI tasks are pre-
 1018 sented in Prompts 1, 2, and 3. The prompt for generating the differences between various reward
 1019 functions is shown in Prompt 4.

1020 Prompt 1: Initial System Prompts of Synthesizing Reward Functions

1021 You are a reward engineer trying to write reward functions to solve reinforcement learning
 1022 tasks as effective as possible.
 1023 Your goal is to write a reward function for the environment that will help the agent learn the
 1024 task described in text.
 1025

1026 Your reward function should use useful variables from the environment as inputs. As an example
 1027 , the reward function signature can be:
 1028 @torch.jit.script
 1029 def compute_reward(object_pos: torch.Tensor, goal_pos: torch.Tensor) -> Tuple[torch.Tensor,
 Dict[str, torch.Tensor]]:
 1030 ...
 1031 return reward, {}
 1032 Since the reward function will be decorated with @torch.jit.script, please make sure that the
 code is compatible with TorchScript (e.g., use torch tensor instead of numpy array).
 1033 Make sure any new tensor or variable you introduce is on the same device as the input tensors.

1034 Prompt 2: Feedback Prompts

1036 The reward function has been iterated {current_iteration} rounds.
 1037 In each iteration, a good reward function and a bad reward function are generated.
 1038 The good reward function generated in the x-th iteration is denoted as "iterx-good", and the
 bad reward function generated is denoted as "iterx-bad".
 1039 The following outlines the differences between these reward functions.

1040 We trained an RL policy using iter1-good reward function code and tracked the values of the
 individual components in the reward function after every {epoch_freq} epochs and the
 maximum, mean, minimum values encountered:
 1041 <REWARD FEEDBACK>
 1042

1043 The difference between iter2-good and iter1-good is: <DIFFERENCE>
 1044 <REPEAT UNTIL THE CURRENT ITERATION>
 1045

1046 Next, the two reward functions generated in the {current_iteration_ordinal} iteration are
 provided.
 1047 The 1st generated reward function is as follows:
 1048 <REWARD FUNCTION>
 1049 We trained an RL policy using the 1st reward function code and tracked the values of the
 individual components in the reward function after every {epoch_freq} epochs and the
 maximum, mean, minimum values encountered:
 1050 <REWARD FEEDBACK>
 1051

1052 The 2nd generated reward function is as follows:
 1053 <REWARD FUNCTION>
 1054 We trained an RL policy using the 2nd reward function code and tracked the values of the
 individual components in the reward function after every {epoch_freq} epochs and the
 maximum, mean, minimum values encountered:
 1055 <REWARD FEEDBACK>
 1056

1057 The following content is the most important information.
 1058 Good example: 1st reward function. Bad example: 2nd reward function.
 You need to modify based on the good example. DO NOT based on the code of the bad example.
 1059 Please carefully analyze the policy feedback and provide a new, improved reward function that
 can better solve the task. Some helpful tips for analyzing the policy feedback:
 1060 (1) If the values for a certain reward component are near identical throughout, then this
 means RL is not able to optimize this component as it is written. You may consider
 1061 (a) Changing its scale or the value of its temperature parameter
 1062 (b) Re-writing the reward component
 1063 (c) Discarding the reward component
 1064 (2) If some reward components' magnitude is significantly larger, then you must re-scale
 its value to a proper range
 1065 Please analyze each existing reward component in the suggested manner above first, and then
 write the reward function code.
 1066

1067 Prompt 3: Prompts of Tips for Writing Reward Functions

1069 The output of the reward function should consist of two items:
 1070 (1) the total reward,
 1071 (2) a dictionary of each individual reward component.
 The code output should be formatted as a python code string: "```python ... ```".

1072 Some helpful tips for writing the reward function code:
 1073 (1) You may find it helpful to normalize the reward to a fixed range by applying
 transformations like torch.exp to the overall reward or its components
 1074 (2) If you choose to transform a reward component, then you must also introduce a
 temperature parameter inside the transformation function; this parameter must be a named
 1075 variable in the reward function and it must not be an input variable. Each transformed
 reward component should have its own temperature variable
 1076 (3) Make sure the type of each input variable is correctly specified; a float input
 variable should not be specified as torch.Tensor
 1077 (4) Most importantly, the reward code's input variables must contain only attributes of
 the provided environment class definition (namely, variables that have prefix self.).
 1078 Under no circumstance can you introduce new input variables.
 1079

Prompt 4: Prompts of Describing Differences

```

1080
1081
1082 You are an engineer skilled at comparing the differences between two reward function code
1083 snippets used in reinforcement learning.
1084 Your goal is to describe the differences between two reward function code snippets.
1085 The following are two reward functions written in Python code used for the task:
1086 <TASK_DESCRIPTION>
1087 The first reward function is as follows:
1088 <REWARD_FUNCTION>
1089 The second reward function is as follows:
1090 <REWARD_FUNCTION>
1091 Please directly describe the differences between these two codes. No additional descriptions
1092 other than the differences are required.

```

D.2 IMO/CPHO/IOI SYSTEM PROMPT

Below we provide the complete system prompt used to guide the Gemini LLM to generate appropriate IMO/CPHO/IOI solutions, perform major vote and choose the best and the worst response.

SYSTEM PROMPT 1

Prompt 5: IMO CoT system prompt

```

1099 -- BEGIN SYSTEM PROMPT --
1100 """
1101 ### Core Instructions ###
1102
1103 * **Rigor is Paramount:** Your primary goal is to produce a complete
1104 and rigorously justified solution. Every step in your solution must
1105 be logically sound and clearly explained. A correct final answer
1106 derived from flawed or incomplete reasoning is considered a failure.
1107 * **Honesty About Completeness:** If you cannot find a complete
1108 solution, you must **not** guess or create a solution that appears
1109 correct but contains hidden flaws or justification gaps. Instead, you
1110 should present only significant partial results that you can
1111 rigorously prove. A partial result is considered significant if it
1112 represents a substantial advancement toward a full solution. Examples
1113 include:
1114 * Proving a key lemma.
1115 * Fully resolving one or more cases within a logically sound case-
1116 based proof.
1117 * Establishing a critical property of the mathematical objects in
1118 the problem.
1119 * For an optimization problem, proving an upper or lower bound
1120 without proving that this bound is achievable.
1121 * **Use TeX for All Mathematics:** All mathematical variables,
1122 expressions, and relations must be enclosed in TeX delimiters (e.g.,
1123 'Let  $n$  be an integer.').
1124
1125 ### Output Format ###
1126
1127 Your response MUST be structured into the following sections, in this
1128 exact order.
1129
1130 *** Final Answer ***
1131
1132 [Your final answer here](You should provide only the final answer here,
1133 without any explanation or reasoning.)
1134
1135 *** Reasoning ***
1136
1137 **1. Summary**
1138
1139 Provide a concise overview of your findings. This section must contain
1140 two parts:

```

```

1134
1135 * **a. Verdict:** State clearly whether you have found a complete
1136 solution or a partial solution.
1137 * **For a complete solution:** State the final answer, e.g., "I
1138 have successfully solved the problem. The final answer is..."
1139 * **For a partial solution:** State the main rigorous conclusion(s)
1140 you were able to prove, e.g., "I have not found a complete solution,
1141 but I have rigorously proven that..."
1142 * **b. Method Sketch:** Present a high-level, conceptual outline of
1143 your solution. This sketch should allow an expert to understand the
1144 logical flow of your argument without reading the full detail. It
1145 should include:
1146 * A narrative of your overall strategy.
1147 * The full and precise mathematical statements of any key lemmas or
1148 major intermediate results.
1149 * If applicable, describe any key constructions or case splits that
1150 form the backbone of your argument.
1151
1152 **2. Detailed Solution**
1153
1154 Present the full, step-by-step mathematical proof. Each step must be
1155 logically justified and clearly explained. The level of detail should
1156 be sufficient for an expert to verify the correctness of your
1157 reasoning without needing to fill in any gaps. This section must
1158 contain ONLY the complete, rigorous proof, free of any internal
1159 commentary, alternative approaches, or failed attempts.
1160
1161 ### Self-Correction Instruction ###
1162
1163 Before finalizing your output, carefully review your "Method Sketch" and
1164 "Detailed Solution" to ensure they are clean, rigorous, and strictly
1165 adhere to all instructions provided above. Verify that every
1166 statement contributes directly to the final, coherent mathematical
1167 argument.
1168
1169 ""
1170
1171 -- END SYSTEM PROMPT --

```

SYSTEM PROMPT 2

Prompt 6: Iterative refinement in 3D Scaling system prompt

```

1172 -- BEGIN SYSTEM PROMPT --
1173
1174 """"You are an expert problem solver.
1175 Your task is to carefully read the problem statement and reflect on two
1176 previous solutions.
1177 - previous_output1 is a relatively better attempt, but it may contain
1178 mistakes or gaps.
1179 - previous_output2 is a weaker attempt, which might include irrelevant
1180 reasoning or errors.
1181
1182 Your job:
1183 1. Identify the strengths and weaknesses of both solutions.
1184 2. Combine the strengths and correct the weaknesses.
1185 3. Produce a new, improved solution that is clearer, more accurate, and
1186 better structured.
1187
1188 Make sure the final answer is complete and stands alone as a polished
1189 solution.""
1190
1191 -- END SYSTEM PROMPT --

```

```

1188
1189
1190 -- BEGIN QUESTION PROMPT --
1191 """
1192 Problem Statement:
1193 {problem_statement}
1194
1195 Better Attempt (previous_output1):
1196 {previous_output1}
1197
1198 Weaker Attempt (previous_output2):
1199 {previous_output2}
1200 """
1201 -- END QUESTION PROMPT --

```

SYSTEM PROMPT 3

Prompt 7: Pairwise Comparison system prompt

```

1206 -- BEGIN SYSTEM PROMPT --
1207
1208 """You are an expert in comparing problem solutions. Your task is to
1209 compare two solutions and output only the better one.
1210
1211 Strictly follow these rules:
1212 1. Only compare the quality of Solution 1 and Solution 2
1213 2. Judge based on accuracy, completeness, and clarity
1214 3. Output must be exactly one of: "solution1" or "solution2"
1215 4. Absolutely do not output any analysis, reasoning, or other text
1216 5. If difficult to judge, choose the relatively more accurate one
1217
1218 Your output must only be: solution1 or solution2"""
1219 -- END SYSTEM PROMPT --
1220
1221 -- BEGIN QUESTION PROMPT --
1222 """Problem statement: {problem_statement}
1223
1224 Solution 1: {result1}
1225
1226 Solution 2: {result2}
1227
1228 Please output the number of the better solution: """
1229 -- END QUESTION PROMPT --

```

SYSTEM PROMPT 4

Prompt 8: Majority Vote system prompt

```

1235 -- BEGIN SYSTEM PROMPT --
1236
1237 """
1238 You are a professional mathematical answer consistency expert. Your task
1239 is to analyze a set of mathematical answers, identify answers that
1240 are essentially the same, and find the most frequently occurring
1241 answer(s) (the mode).
1242
1243 # Core Principles

```

```

1242 The criterion for judging whether two answers are the same is: whether
1243 they are mathematically equivalent, not whether the strings are
1244 exactly the same.
1245
1246 # Equivalence Rules
1247 1. **Numerical equivalence**:  $0.75 = 3/4 = 75\% = \frac{3}{4} = \text{"three quarters"}$ 
1248 2. **Algebraic expression equivalence**:  $2x + 3 = 3 + 2x = (4x + 6)/2$ 
1249 3. **Set equivalence**:  $\{1, 2, 3\} = \{3, 2, 1\} = \{x \mid x \in \{1,2,3\}\}$ 
1250 4. **Interval equivalence**:  $(0,1) = \{x \mid 0 < x < 1\} = \text{"open interval from 0 to 1"}$ 
1251 5. **Function equivalence**:  $f(x) = x^2 = x*x = x^2$ 
1252 6. **Geometric equivalence**:  $\text{"right triangle"} = \text{"triangle with a 90 degree angle"}$ 
1253 7. **Logical equivalence**:  $\text{true} = \text{"correct"}$ 
1254
1255 # Handling natural language answers
1256 For answers containing explanations, extract the core mathematical
1257 content:
1258 -  $\text{"The answer is } 3/4 \text{ because..."} \rightarrow \text{extract } \frac{3}{4}$ 
1259 -  $\text{"I think it should be } 2\pi \text{"} \rightarrow \text{extract } 2\pi$ 
1260 -  $\text{"The area of this triangle is 12 square centimeters"} \rightarrow \text{extract } 12$ 
1261
1262 # Output requirements
1263 1. **Return only the mode answer(s)**, no explanation
1264 2. **Return in the most concise standard form** (prefer mathematical
1265 symbols)
1266 3. **If there are multiple modes** (same highest frequency), separate
1267 them with commas
1268 4. **Keep original format**: if it's a set, return in set form; if
1269 interval, return interval form
1270
1271 # Examples
1272 Input: ["0.75", "3/4", "75%", "The answer is three quarters"]
1273 Output:  $\frac{3}{4}$ 
1274
1275 Input: ["{1,2,3}", "{3,1,2}", "set contains 1,2,3"]
1276 Output: {1,2,3}
1277
1278 Input: ["(0,\infty)", "x>0", "positive real numbers"]
1279 Output:  $(0, \infty)$ 
1280
1281 Input: ["2", "2.0", "two", "The answer is 2"]
1282 Output: 2
1283 ""
1284
1285 -- END SYSTEM PROMPT --

```

SYSTEM PROMPT 5

Prompt 9: CPHO CoT system prompt

```

1288 -- BEGIN SYSTEM PROMPT --
1289
1290 ""
1291 You are a professional physicist with expertise in solving high school
1292 and undergraduate level physics problems. Your task is to provide a
1293 complete, rigorous, and well-justified solution to the given physics
1294 problem.
1295
1296 ### Core Instructions ###

```

1296 * ****Complete Coverage is Paramount:**** Your primary goal is to produce a
1297 complete and rigorously justified solution for every sub-question (
1298 each marked with ‘\item’ in the problem statement). You must answer
1299 all sub-questions in the order they are presented. Do not skip any
1300 sub-question or terminate early after answering only a subset. Each
1301 sub-question’s solution must be logically sound, physically accurate,
1302 and clearly explained.

1303 * ****Rigor and Detail:**** For each sub-question, provide a step-by-step
1304 detailed process that includes all reasoning, calculations, and
1305 physical principles applied. All mathematical variables, expressions,
1306 and relations must be enclosed in TeX delimiters (e.g., ‘ $F = ma$ ’).
1307 Ensure that units, dimensions, and significant figures are handled
1308 appropriately where relevant.

1309 * ****Honesty About Completeness:**** If you cannot solve a sub-question
1310 completely, you must not guess or create an answer that appears
1311 correct but contains flaws. Instead, present any partial results you
1312 can rigorously justify, and clearly indicate which sub-question
1313 remains unsolved or partially solved. A partial result should
1314 represent a substantial advancement, such as deriving a key equation
1315 or setting up a correct problem framework.

1316 * ****Final Answers Listing:**** After completing the detailed solutions
1317 for all sub-questions, you must list all final answers in order at
1318 the very end of your response. This listing should include only the
1319 answers (e.g., numerical values, expressions, or conclusions),
1320 without the detailed processes.

1321 * ****Please notice:**** If there is a sub-question marked as "key sub-
1322 question", the final answer to that sub-question should be
1323 highlighted as the "Key Final Answer" in your final answers listing.
1324 If there is not such a sub-question, please treat the last sub-
1325 question as the key one. Your output should follow the structure
1326 below.

1327 **### Output Format ###**

1328 Your response MUST be structured into the following sections, in this
1329 exact order.

1330 ***** Key Final Answer *****
1331 [The Key Final Answer]
1332 (In this section, provide the final answer of the key sub-question only.
1333 In the problem statement part, there would be a sub-question marked
1334 as "key sub-question".
1335 If there is no such sub-question, please list the final answer of the
1336 last sub-question in this section.)

1337 ***** All Final Answers *****
1338 List all final answers in order, corresponding to each sub-question. This
1339 section should be concise and contain only the answers, formatted as
1340 :

1341 * Sub-question 1: [Answer]
1342 * Sub-question 2: [Answer]
1343 * ... and so on for all sub-questions.

1344 ***** Reasoning *****

1345 Present the full, step-by-step solutions for each sub-question in
1346 sequence. For each sub-question:

1347 * Start with a clear heading indicating the sub-question number or
1348 label (e.g., "Sub-question 1:").
1349 * Provide a rigorous and detailed solution, including all reasoning,
1350 calculations, and explanations. Use TeX for mathematics.

```

1350 * Ensure that each step is justified physically and mathematically. If
1351 a sub-question builds on previous answers, reference them
1352 appropriately.
1353 * Do not include commentary on alternative approaches or failed
1354 attempts-only the coherent argument for each sub-question.
1355
1356 ### Self-Correction Instruction ###
1357
1358 Before finalizing your output, carefully review your response to ensure:
1359 - All sub-questions have been addressed in the order presented, with no
1360 omissions.
1361 - Each detailed solution is complete, rigorous, and free of gaps.
1362 - The final answers are accurately derived and listed correctly at the
1363 end.
1364 - The output adheres strictly to this format and instructions.
1365
1366 -- END SYSTEM PROMPT --

```

SYSTEM PROMPT 6

Prompt 10: IOI CoT system prompt

```

1370 -- BEGIN SYSTEM PROMPT --
1371
1372 ""
1373
1374 ### Core Instructions ###
1375
1376 * **Rigor is Paramount:** Your primary goal is to produce a **fully
1377 correct and executable** C++ code. The code must handle all valid
1378 inputs defined in the problem statement and must explicitly deal with
1379 edge cases. You should also provide a detailed explanation of your
1380 algorithm in your code to demonstrate your main method and why it is
1381 correct.
1382 * **Honesty About Completeness:** If you cannot provide a complete,
1383 correct code implementation, you must not guess or conceal flaws.
1384 Instead, present only the significant partial results that you can
1385 rigorously justify. For example:
1386 - A code that can solve subtasks with the highest total score, you
1387 should make sure its correct and provide its main algorithm.
1388 - A possible algorithm direction that can solve the whole problem
1389 although you do not implement it correctly.
1390 - A correct implementation of a critical function or subroutine.
1391 * **Rule for Function Call:** If the problem involves invoking
1392 functions that you are not required to implement, you must ensure
1393 that every invocation strictly adheres to the problems
1394 specifications; otherwise, your code will be deemed invalid. Each
1395 invocation may alter the state of the data in ways that affect your
1396 objectives, and once made, such calls cannot be undone
1397 * **Use TeX for All Mathematics:** All mathematical variables,
1398 expressions, and relations in your algorithm must be enclosed in TeX
1399 delimiters (e.g., 'Let  $n$  be an integer.').
1400 * **Code Format:** Your code should read the inputs from stdin solve
1401 the problem and write the answer to stdout (do not directly test on
1402 the sample inputs). Enclose your code within delimiters as follows.
1403 Ensure your c++ program contains the function required in the
1404 problem statement.\n``cpp\n// YOUR CODE HERE\n``"
1405
1406 ### Output Format ###
1407
1408 Your response MUST be structured into the following sections, in this
1409 exact order.

```

```

1404 **1. Summary**
1405
1406 Provide a concise overview of your findings. This section must contain
1407 two parts:
1408
1409 * **a. Verdict:** State clearly whether you have found a complete
1410 solution or a partial solution.
1411 * **For a complete solution:** State the final code, e.g., "I have
1412 successfully solved the problem. The final code is ..."
1413 * **For a partial solution:** State the partial code you now have,
1414 e.g., "I have not found a complete solution, but I have a code that
1415 can solve subtasks with the highest total score, the code is ```cpp
1416 ... ```"
1417 * **b. Method Sketch:** Present a high-level, conceptual outline of
1418 your algorithm. This sketch should allow an expert to understand the
1419 main algorithm of your argument without reading the full detail.
1420
1421 **2. Detailed Solution**
1422
1423 Present the full, step-by-step explanation of your code.
1424 If your algorithm requires some proof on complexity or correctness, you
1425 should also provide the proof.
1426 If your answer contains algorithms that can solve subtasks, you should
1427 also describe them.
1428 The level of detail should be sufficient for an expert to verify the
1429 correctness of your code without needing to test it in testcase.
1430
1431 **3. Final Code**
1432
1433 Present your final code for the problem again. Place the solution inside
1434 one fenced code block (### Answer: (use the provided format with
1435 backticks)```cpp ...```).
1436
1437 ### Self-Correction Instruction ###
1438
1439 Before finalizing your output, carefully review your code and algorithm.
1440 Fix any bugs, make sure the code is executable.
1441
1442 -- END SYSTEM PROMPT --

```

SYSTEM PROMPT 7

Prompt 11: IOI CoT system prompt

```

1443 -- BEGIN SYSTEM PROMPT --
1444
1445 """
1446 You are an expert in evaluating problem solutions. Your task is to select
1447 the single best solution among several options. Strictly follow
1448 these rules:
1449 1. Compare all provided solutions based on accuracy, completeness,
1450 clarity, and overall quality. 2. Output only the number of the best
1451 solution (starting from 1). 3. Do not output any reasoning,
1452 explanations, or extra text. 4. If it is difficult to decide, choose
1453 the solution that is relatively more accurate and complete.
1454 output format:
1455 "Solution 1" or "Solution 2" or ... (just output "Solution" and one
1456 number following)
1457 Your output must be exactly the number of the best solution.
1458 """
1459
1460 -- END SYSTEM PROMPT --

```

```

1458 -- BEGIN QUESTION PROMPT --
1459
1460 solutions_text = "\n\n".join([f"Solution {i+1}: {s}" for i, s in
1461     enumerate(results)])
1462
1463 question_prompt = f"""
1464 Problem statement: {problem_statement}
1465
1466 {solutions_text}
1467
1468 Please output only the number of the best solution (starting from 1):"""
1469 -- END QUESTION PROMPT --

```

1470

1471 D.3 IOI SYSTEM PROMPT

1472

1473 Below we provide the complete system prompt used to guide the Gemini LLM to generate appropriate solutions for IOI problem.

1474

1476 D.4 HUMAN-IN-THE-LOOP PREFERENCE

1477

1478 D.4.1 DEMOGRAPHIC DATA

1479

1480 The participants in the human-in-the-loop preference experiments consisted of 7 individuals aged 19 to 30, including 2 women and 5 men. Their educational backgrounds included 2 undergraduate students and 5 graduate students. The 20 volunteers recruited to evaluate the performance of different methods were aged 23 to 28, comprising 5 women and 15 men, with 3 undergraduates and 17 graduate students.

1481

1485 D.4.2 ISAACGYM TASKS

1486

1487 We evaluate human-in-the-loop preference experiments on tasks in IsaacGym, including *Quadcopter*, *Humanoid*, *Ant*, *ShadowHand*, and *AllegroHand*. In these experiments, volunteers were limited to comparing reward functions based solely on videos showcasing the final policies derived from each reward function.

1488

1489 In the *Quadcopter* task, humans evaluate performance by observing whether the quadcopter moves quickly and efficiently, and whether it stabilizes in the final position. For the *Humanoid* and *Ant* tasks, where the task description is "make the ant/humanoid run as fast as possible," humans estimate speed by comparing the time taken to cover the same distance and assessing the movement posture. However, due to the variability in movement postures and directions, estimating speed can introduce inaccuracies. In the *ShadowHand* and *AllegroHand* tasks, where the goal is "to make the hand spin the object to a target orientation," Humans find it challenging to calculate the precise difference between the current orientation and the target orientation at every moment, even though the target orientation is displayed nearby. Nevertheless, humans still can estimate the duration of effective rotations with the target orientation in the video, thus evaluating the performance of a single spin. Since the target orientation regenerates upon being reached, the frequency of target orientation changes can also aid in facilitating the assessment of evaluating performance.

1490

1491 Due to the lack of precise environmental data, volunteers cannot make absolutely accurate judgments during the experiments. For instance, in the *Humanoid* task, robots may move in varying directions, which can introduce biases in volunteers' assessments of speed. However, volunteers are still able to filter out extremely poor results and select videos with relatively better performance. In most cases, the selected results closely align with those derived from proxy human preferences, enabling effective improvements in task performance.

1492

1493 Below is a specific case from the *Humanoid* task that illustrates the potential errors humans may make during evaluation and the learning process of the reward function under this assumption. The reward task scores (RTS) chosen by the volunteer across five iterations are 4.521, 6.069, 6.814, 6.363, 6.983.

1494

1512 In the first iteration, the ground-truth task scores of each policy were
 1513 0.593, 2.744, 4.520, 0.192, 2.517, 5.937, although the volunteer was unaware of these scores.
 1514 Initially, the volunteer eliminated policies 0 and 3, as the robots in those videos primarily exhibited
 1515 spinning behavior. Subsequently, the volunteer assessed the speed of the remaining robots based
 1516 on how quickly a specific robot moved out of the field. The volunteer correctly identified that the
 1517 robots in policies 1 and 4 were slightly slower. However, due to minor differences in the movement
 1518 directions of the robots in policies 2 and 5, the volunteer mistakenly selected policy 2 as the best
 1519 option, incorrectly concluding that the robot in policy 2 was faster.

1520 Thus, the reward function selected in iteration 1 consists of several key components: velocity reward,
 1521 upright reward, force penalty, unnatural pose penalty, and action penalty. These components not only
 1522 promote faster training, which is the primary objective, but also encourage the maintenance of an
 1523 upright pose. Additionally, the function penalizes excessive force usage, extreme joint angles, and
 1524 large action values to foster smoother and more controlled movements.

1525 In subsequent iterations, the volunteer effectively identified reward functions that exhibited relatively
 1526 better and worse performance outcomes. Adjustments were made to the weights of each component,
 1527 and specific temperature values were introduced for each. These modifications resulted in a more
 1528 balanced reward structure, ensuring that critical aspects exert a stronger influence, thereby allowing
 1529 for greater control over the learning dynamics and improving the agent’s performance in achieving
 1530 the task. Even in Iteration 4, the volunteer did not select the reward function with the highest RTS
 1531 (6.813) but instead opted for the second-highest reward function (RTS = 6.363). Nevertheless, the
 1532 reward function exhibited consistent improvement during these iterations.

1533 Here we show the full reward function during the process.

1534

Humanoid Task: Reward Function chosen by volunteer with RTS (4.521) of Iteration 1

```

1537 def compute_reward(
1538     velocity: torch.Tensor,
1539     dof_pos: torch.Tensor,
1540     dof_force_tensor: torch.Tensor,
1541     actions: torch.Tensor,
1542     up_vec: torch.Tensor,
1543     heading_vec: torch.Tensor
1544 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
1545     # Reward for forward velocity (run as fast as possible)
1546     velocity_reward = velocity[:, 0]
1547
1548     # Encouragement for upright pose (penalize for deviation from
1549     # vertical)
1550     upright_reward = up_vec[:, 2]
1551
1552     # Penalize high force usage (energy efficiency)
1553     force_penalty = torch.sum(torch.abs(dof_force_tensor), dim=1)
1554
1555     # Penalize unnatural joint positions (for instance, avoid
1556     # extreme angles)
1557     unnatural_pose_penalty = torch.sum(torch.abs(dof_pos), dim=1)
1558
1559     # Penalize for large actions (to promote smoother movement)
1560     action_penalty = torch.sum(torch.abs(actions), dim=1)
1561
1562     # Normalize the rewards and penalties
1563     velocity_reward = torch.exp(velocity_reward) - 1
1564     upright_reward = torch.exp(upright_reward) - 1
1565     temperature = 1.0
1566     force_penalty = torch.exp(-force_penalty / temperature)
1567     unnatural_pose_penalty = torch.exp(-unnatural_pose_penalty /
1568     temperature)
1569     action_penalty = torch.exp(-action_penalty / temperature)
  
```

```

1566
1567 # Combine the rewards and penalties into a single reward
1568 total_reward = (
1569     velocity_reward +
1570     0.5 * upright_reward -
1571     0.01 * force_penalty -
1572     0.01 * unnatural_pose_penalty -
1573     0.01 * action_penalty
1574 )
1575 # Return the total reward and each component for analysis
1576 reward_components = {
1577     "velocity_reward": velocity_reward,
1578     "upright_reward": upright_reward,
1579     "force_penalty": force_penalty,
1580     "unnatural_pose_penalty": unnatural_pose_penalty,
1581     "action_penalty": action_penalty
1582 }
1583
1584 return total_reward, reward_components

```

Humanoid Task: Reward Function chosen by volunteer with RTS (6.069) of Iteration 2

```

1587
1588 def compute_reward(
1589     velocity: torch.Tensor,
1590     dof_pos: torch.Tensor,
1591     dof_force_tensor: torch.Tensor,
1592     actions: torch.Tensor,
1593     up_vec: torch.Tensor,
1594     heading_vec: torch.Tensor
1595 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
1596     # Reward for forward velocity (run as fast as possible)
1597     velocity_reward = velocity[:, 0]
1598     velocity_temperature = 1.2 # increased slightly
1599     velocity_reward = torch.exp(velocity_reward /
1600     velocity_temperature) - 1
1601
1602     # Encouragement for upright pose (penalize for deviation from
1603     vertical)
1604     upright_reward = up_vec[:, 2]
1605     upright_temperature = 0.5 # introduce a specific temperature
1606     upright_reward = torch.exp(upright_reward / upright_temperature
1607     ) - 1
1608
1609     # Penalize high force usage (energy efficiency)
1610     force_penalty = torch.sum(torch.abs(dof_force_tensor), dim=1)
1611     force_temperature = 0.1 # decreased to make it more
1612     significant
1613     force_penalty = torch.exp(-force_penalty / force_temperature)
1614
1615     # Penalize unnatural joint positions (for instance, avoid
1616     extreme angles)
1617     unnatural_pose_penalty = torch.sum(torch.abs(dof_pos), dim=1)
1618     pose_temperature = 0.1 # decreased to make it more significant
1619     unnatural_pose_penalty = torch.exp(-unnatural_pose_penalty /
1620     pose_temperature)
1621
1622     # Penalize for large actions (to promote smoother movement)
1623     action_penalty = torch.sum(torch.abs(actions), dim=1)
1624     action_temperature = 0.1 # decreased to make it more
1625     significant

```

```

1620
1621     action_penalty = torch.exp(-action_penalty / action_temperature
1622 )
1623
1624     # Combine the rewards and penalties into a single reward
1625     total_reward = (
1626         velocity_reward +
1627         0.5 * upright_reward -
1628         0.01 * force_penalty -
1629         0.01 * unnatural_pose_penalty -
1630     )
1631
1632     # Return the total reward and each component for analysis
1633     reward_components = {
1634         "velocity_reward": velocity_reward,
1635         "upright_reward": upright_reward,
1636         "force_penalty": force_penalty,
1637         "unnatural_pose_penalty": unnatural_pose_penalty,
1638         "action_penalty": action_penalty
1639     }
1640
1641     return total_reward, reward_components

```

Humanoid Task: Reward Function chosen by volunteer with RTS (6.814) of Iteration 3

```

1644 def compute_reward(
1645     velocity: torch.Tensor,
1646     dof_pos: torch.Tensor,
1647     dof_force_tensor: torch.Tensor,
1648     actions: torch.Tensor,
1649     up_vec: torch.Tensor,
1650     heading_vec: torch.Tensor
1651 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
1652     # Reward for forward velocity (run as fast as possible)
1653     velocity_reward = velocity[:, 0]
1654     velocity_temperature = 1.1 # minor adjustment
1655     velocity_reward = torch.exp(velocity_reward /
1656         velocity_temperature) - 1
1657
1658     # Encouragement for upright pose (penalize for deviation from
1659     # vertical)
1660     upright_reward = up_vec[:, 2]
1661     upright_temperature = 0.6 # slight adjustment
1662     upright_reward = torch.exp(upright_reward / upright_temperature)
1663     - 1
1664
1665     # Penalize high force usage (energy efficiency)
1666     force_penalty = torch.sum(torch.abs(dof_force_tensor), dim=1)
1667     force_temperature = 0.15 # increased to try to make it
1668     # effective
1669     force_penalty = torch.exp(-force_penalty / force_temperature)
1670
1671     # Penalize unnatural joint positions (for instance, avoid
1672     # extreme angles)
1673     unnatural_pose_penalty = torch.sum(torch.abs(dof_pos), dim=1)
1674     pose_temperature = 0.2 # increased to try to make it effective
1675     unnatural_pose_penalty = torch.exp(-unnatural_pose_penalty /
1676         pose_temperature)
1677
1678     # Penalize for large actions (to promote smoother movement)

```

```

1674
1675 action_penalty = torch.sum(torch.abs(actions), dim=1)
1676 action_temperature = 0.2 # increased to try to make it
1677 effective
1678 action_penalty = torch.exp(-action_penalty / action_temperature
1679 )
1680
1681 # Combine the rewards and penalties into a single reward
1682 total_reward = (
1683     velocity_reward +
1684     0.5 * upright_reward -
1685     0.02 * force_penalty - # increased slightly for more
1686     impact
1687     0.02 * unnatural_pose_penalty - # increased slightly for
1688     more impact
1689     0.02 * action_penalty # increased slightly for more impact
1690 )
1691
1692 # Return the total reward and each component for analysis
1693 reward_components = {
1694     "velocity_reward": velocity_reward,
1695     "upright_reward": upright_reward,
1696     "force_penalty": force_penalty,
1697     "unnatural_pose_penalty": unnatural_pose_penalty,
1698     "action_penalty": action_penalty
1699 }
1700
1701 return total_reward, reward_components

```

Humanoid Task: Reward Function chosen by volunteer with RTS (6.363) of Iteration 4

```

1703
1704 def compute_reward(
1705     velocity: torch.Tensor,
1706     dof_pos: torch.Tensor,
1707     dof_force_tensor: torch.Tensor,
1708     actions: torch.Tensor,
1709     up_vec: torch.Tensor,
1710     heading_vec: torch.Tensor
1711 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
1712     # Reward for forward velocity (run as fast as possible)
1713     velocity_reward = velocity[:, 0]
1714     velocity_temperature = 1.05 # slight adjustment to refine the
1715     impact
1716     velocity_reward = torch.exp(velocity_reward /
1717     velocity_temperature) - 1
1718
1719     # Encouragement for upright pose (penalize for deviation from
1720     vertical)
1721     upright_reward = up_vec[:, 2]
1722     upright_temperature = 0.65 # slight loosening for more upright
1723     reward
1724     upright_reward = torch.exp(upright_reward / upright_temperature
1725     ) - 1
1726
1727     # Penalize high force usage (energy efficiency)
1728     force_penalty = torch.sum(torch.abs(dof_force_tensor), dim=1)
1729     force_temperature = 0.2 # increased to make it more
1730     significant
1731     force_penalty = torch.exp(-force_penalty / force_temperature)

```

```

1728
1729 # Penalize unnatural joint positions (for instance, avoid
1730 extreme angles)
1731 unnatural_pose_penalty = torch.sum(torch.abs(dof_pos), dim=1)
1732 pose_temperature = 0.25 # slight increase to make this
1733 component effective
1734 unnatural_pose_penalty = torch.exp(-unnatural_pose_penalty /
1735 pose_temperature)
1736
1737 # Penalize for large actions (to promote smoother movement)
1738 action_penalty = torch.sum(torch.abs(actions), dim=1)
1739 action_temperature = 0.25 # slightly adjusted for more
1740 prominent constraint
1741 action_penalty = torch.exp(-action_penalty / action_temperature
1742 )
1743
1744 # Combine the rewards and penalties into a single reward
1745 total_reward = (
1746     velocity_reward +
1747     0.5 * upright_reward -
1748     0.015 * force_penalty - # slight increase for more impact
1749     0.015 * unnatural_pose_penalty - # slight increase for
1750 more impact
1751     0.015 * action_penalty # slight increase for more impact
1752 )
1753
1754 # Return the total reward and each component for analysis
1755 reward_components = {
1756     "velocity_reward": velocity_reward,
1757     "upright_reward": upright_reward,
1758     "force_penalty": force_penalty,
1759     "unnatural_pose_penalty": unnatural_pose_penalty,
1760     "action_penalty": action_penalty
1761 }
1762
1763 return total_reward, reward_components

```

Humanoid Task: Reward Function with best RTS (6.813) of Iteration 4(not chosen by volunteer)

```

1764 def compute_reward(
1765     velocity: torch.Tensor,
1766     dof_pos: torch.Tensor,
1767     dof_force_tensor: torch.Tensor,
1768     actions: torch.Tensor,
1769     up_vec: torch.Tensor,
1770     heading_vec: torch.Tensor
1771 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
1772     # Reward for forward velocity (run as fast as possible)
1773     velocity_reward = velocity[:, 0]
1774     velocity_temperature = 1.15
1775     velocity_reward = torch.exp(velocity_reward /
1776     velocity_temperature) - 1
1777
1778     # Encouragement for upright pose (penalize for deviation from
1779     vertical)
1780     upright_reward = up_vec[:, 2]
1781     upright_temperature = 0.55
1782     upright_reward = torch.exp(upright_reward / upright_temperature
1783     ) - 1

```

```

1782
1783 # Penalize high force usage (energy efficiency)
1784 force_penalty = torch.sum(torch.abs(dof_force_tensor), dim=1)
1785 force_temperature = 0.12
1786 force_penalty = torch.exp(-force_penalty / force_temperature)
1787
1788 # Penalize unnatural joint positions (for instance, avoid
1789 extreme angles)
1790 unnatural_pose_penalty = torch.sum(torch.abs(dof_pos), dim=1)
1791 pose_temperature = 0.18
1792 unnatural_pose_penalty = torch.exp(-unnatural_pose_penalty /
1793 pose_temperature)
1794
1795 # Penalize for large actions (to promote smoother movement)
1796 action_penalty = torch.sum(torch.abs(actions), dim=1)
1797 action_temperature = 0.18
1798 action_penalty = torch.exp(-action_penalty / action_temperature
1799 )
1800
1801 # Combine the rewards and penalties into a single reward
1802 total_reward = (
1803     velocity_reward +
1804     0.5 * upright_reward -
1805     0.02 * force_penalty -
1806     0.02 * unnatural_pose_penalty -
1807     0.02 * action_penalty
1808 )
1809
1810 # Return the total reward and each component for analysis
1811 reward_components = {
1812     "velocity_reward": velocity_reward,
1813     "upright_reward": upright_reward,
1814     "force_penalty": force_penalty,
1815     "unnatural_pose_penalty": unnatural_pose_penalty,
1816     "action_penalty": action_penalty
1817 }
1818
1819 return total_reward, reward_components

```

Humanoid Task: Reward Function chosen by volunteer with RTS (6.983) of Iteration 5

```

1819 def compute_reward(
1820     velocity: torch.Tensor,
1821     dof_pos: torch.Tensor,
1822     dof_force_tensor: torch.Tensor,
1823     actions: torch.Tensor,
1824     up_vec: torch.Tensor,
1825     heading_vec: torch.Tensor
1826 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
1827     # Adjusted parameters based on analysis
1828     velocity_temperature = 1.1
1829     upright_temperature = 0.65
1830     force_temperature = 0.25
1831     pose_temperature = 0.3
1832     action_temperature = 0.3
1833
1834     # Reward for forward velocity (run as fast as possible)
1835     velocity_reward = velocity[:, 0]
1836     velocity_reward = torch.exp(velocity_reward /
1837     velocity_temperature) - 1

```

```

1836
1837 # Encouragement for upright pose (penalize for deviation from
1838 vertical)
1839 upright_reward = up_vec[:, 2]
1840 upright_reward = torch.exp(upright_reward / upright_temperature
1841 ) - 1
1842
1843 # Penalize high force usage (energy efficiency)
1844 force_penalty = torch.sum(torch.abs(dof_force_tensor), dim=1)
1845 force_penalty = torch.exp(-force_penalty / force_temperature)
1846
1847 # Penalize unnatural joint positions (for instance, avoid
1848 extreme angles)
1849 unnatural_pose_penalty = torch.sum(torch.abs(dof_pos), dim=1)
1850 unnatural_pose_penalty = torch.exp(-unnatural_pose_penalty /
1851 pose_temperature)
1852
1853 # Penalize for large actions (to promote smoother movement)
1854 action_penalty = torch.sum(torch.abs(actions), dim=1)
1855 action_penalty = torch.exp(-action_penalty / action_temperature
1856 )
1857
1858 # Combine the rewards and penalties into a single reward
1859 total_reward = (
1860     velocity_reward +
1861     0.5 * upright_reward -
1862     0.02 * force_penalty -
1863     0.02 * unnatural_pose_penalty -
1864     0.02 * action_penalty
1865 )
1866
1867 # Return the total reward and each component for analysis
1868 reward_components = {
1869     "velocity_reward": velocity_reward,
1870     "upright_reward": upright_reward,
1871     "force_penalty": force_penalty,
1872     "unnatural_pose_penalty": unnatural_pose_penalty,
1873     "action_penalty": action_penalty
1874 }
1875
1876 return total_reward, reward_components
1877

```

1873 D.4.3 HUMANOIDJUMP TASK

1874 In our study, we introduced a novel task: *HumanoidJump*, with the task description being “to make
1875 humanoid jump like a real human.” The prompt of environment context in this task is shown in
1876 Prompt 12.

1878 Prompt 12: Prompts of Environment Context in *HumanoidJump* Task

```

1879 class HumanoidJump(VecTask):
1880     """Rest of the environment definition omitted."""
1881     def compute_observations(self):
1882         self.gym.refresh_dof_state_tensor(self.sim)
1883         self.gym.refresh_actor_root_state_tensor(self.sim)
1884         self.gym.refresh_force_sensor_tensor(self.sim)
1885         self.gym.refresh_dof_force_tensor(self.sim)
1886
1887         self.obs_buf[:, self.torso_position:],
1888         self.prev_torso_position[:, self.velocity_world:],
1889         self.angular_velocity_world[:, self.velocity_local:],
1890         self.angular_velocity_local[:, self.up_vec:],
1891         self.heading_vec[:, self.right_leg_contact_force:],
1892         self.left_leg_contact_force[:] = \
1893             compute_humanoid_jump_observations(
1894                 self.obs_buf, self.root_states, self.torso_position,
1895                 self.inv_start_rot, self.dof_pos, self.dof_vel,

```

```

1890         self.dof_force_tensor, self.dof_limits_lower,
1891         self.dof_limits_upper, self.dof_vel_scale,
1892         self.vec_sensor_tensor, self.actions,
1893         self.dt, self.contact_force_scale,
1894         self.angular_velocity_scale,
1895         self.basis_vec0, self.basis_vec1)
1896
1897     def compute_humanoid_jump_observations(obs_buf, root_states, torso_position, inv_start_rot
1898     , dof_pos, dof_vel, dof_force, dof_limits_lower, dof_limits_upper, dof_vel_scale,
1899     sensor_force_torques, actions, dt, contact_force_scale, angular_velocity_scale,
1900     basis_vec0, basis_vec1):
1901         # type: (Tensor, Tensor, Tensor, Tensor, Tensor, Tensor, Tensor, Tensor, Tensor, float
1902         , Tensor, Tensor, float, float, float, Tensor, Tensor) -> Tuple[Tensor, Tensor, Tensor,
1903         Tensor, Tensor, Tensor, Tensor, Tensor, Tensor, Tensor]
1904
1905         prev_torso_position_new = torso_position.clone()
1906
1907         torso_position = root_states[:, 0:3]
1908         torso_rotation = root_states[:, 3:7]
1909         velocity_world = root_states[:, 7:10]
1910         angular_velocity_world = root_states[:, 10:13]
1911
1912         torso_quat, up_proj, up_vec, heading_vec = compute_heading_and_up_vec(
1913             torso_rotation, inv_start_rot, basis_vec0, basis_vec1, 2)
1914
1915         velocity_local, angular_velocity_local, roll, pitch, yaw = compute_rot_new(
1916             torso_quat, velocity_world, angular_velocity_world)
1917
1918         roll = normalize_angle(roll).unsqueeze(-1)
1919         yaw = normalize_angle(yaw).unsqueeze(-1)
1920         dof_pos_scaled = unscale(dof_pos, dof_limits_lower, dof_limits_upper)
1921         scale_angular_velocity_local = angular_velocity_local * angular_velocity_scale
1922
1923         obs = torch.cat((root_states[:, 0:3].view(-1, 3), velocity_local,
1924             scale_angular_velocity_local,
1925             yaw, roll, up_proj.unsqueeze(-1),
1926             dof_pos_scaled, dof_vel * dof_vel_scale,
1927             dof_force * contact_force_scale,
1928             sensor_force_torques.view(-1, 12) * contact_force_scale,
1929             actions), dim=-1)
1930
1931         right_leg_contact_force = sensor_force_torques[:, 0:3]
1932         left_leg_contact_force = sensor_force_torques[:, 6:9]
1933
1934         abdomen_y_pos = dof_pos[:, 0]
1935         abdomen_z_pos = dof_pos[:, 1]
1936         abdomen_x_pos = dof_pos[:, 2]
1937         right_hip_x_pos = dof_pos[:, 3]
1938         right_hip_z_pos = dof_pos[:, 4]
1939         right_hip_y_pos = dof_pos[:, 5]
1940         right_knee_pos = dof_pos[:, 6]
1941         right_ankle_x_pos = dof_pos[:, 7]
1942         right_ankle_y_pos = dof_pos[:, 8]
1943         left_hip_x_pos = dof_pos[:, 9]
1944         left_hip_z_pos = dof_pos[:, 10]
1945         left_hip_y_pos = dof_pos[:, 11]
1946         left_knee_pos = dof_pos[:, 12]
1947         left_ankle_x_pos = dof_pos[:, 13]
1948         left_ankle_y_pos = dof_pos[:, 14]
1949         right_shoulder1_pos = dof_pos[:, 15]
1950         right_shoulder2_pos = dof_pos[:, 16]
1951         right_elbow_pos = dof_pos[:, 17]
1952         left_shoulder1_pos = dof_pos[:, 18]
1953         left_shoulder2_pos = dof_pos[:, 19]
1954         left_elbow_pos = dof_pos[:, 20]
1955
1956         right_shoulder1_action = actions[:, 15]
1957         right_shoulder2_action = actions[:, 16]
1958         right_elbow_action = actions[:, 17]
1959         left_shoulder1_action = actions[:, 18]
1960         left_shoulder2_action = actions[:, 19]
1961         left_elbow_action = actions[:, 20]
1962
1963         return obs, torso_position, prev_torso_position_new, velocity_world,
1964             angular_velocity_world, velocity_local, scale_angular_velocity_local,
1965             up_vec, heading_vec, right_leg_contact_force, left_leg_contact_force

```

1942 **Reward functions.** We show the reward functions in a trial that successfully evolved a human-like
1943 jump: bending both legs to jump. Initially, the reward function focused on encouraging vertical

1944 movement while penalizing horizontal displacement, high contact force usage, and improper joint
 1945 movements. Over time, the scaling factors for the rewards and penalties were gradually adjusted
 1946 by changing the temperature parameters in the exponential scaling. These adjustments aimed to en-
 1947 hance the model’s sensitivity to different movement behaviors. For example, the vertical movement
 1948 reward’s temperature was reduced, leading to more precise rewards for positive vertical movements.
 1949 Similarly, the horizontal displacement penalty was fine-tuned by modifying its temperature across
 1950 iterations, either decreasing or increasing the penalty’s impact on lateral movements. The contact
 1951 force penalty evolved by decreasing its temperature to penalize excessive force usage more strongly,
 1952 especially in the later iterations, making the task more sensitive to leg contact forces. Finally, the
 1953 joint usage reward was refined by adjusting the temperature to either encourage or discourage cer-
 1954 tain joint behaviors, with more focus on leg extension and contraction patterns. Overall, the changes
 1955 primarily revolved around adjusting the sensitivity of different components, refining the balance
 1956 between rewards and penalties to better align the humanoid’s behavior with the desired jumping
 1957 performance.

HumanoidJump Task: Reward Function of Iteration 1

```

1960 def compute_reward(torso_position: torch.Tensor,
1961                   prev_torso_position: torch.Tensor, velocity_world: torch.Tensor,
1962                   right_leg_contact_force: torch.Tensor,
1963                   left_leg_contact_force: torch.Tensor, dof_pos: torch.Tensor) ->
1964   Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
1965     # Ensure all tensors are on the same device
1966     device = torso_position.device
1967
1968     # Compute vertical torso movement reward
1969     vertical_movement = torso_position[:, 2] - prev_torso_position
1970     vertical_movement_reward = torch.clamp(vertical_movement, min
1971     =0.0) # Reward positive vertical movement
1972     vertical_movement_reward = torch.exp(vertical_movement_reward /
1973     0.1) # Use exponential scaling with temperature
1974
1975     # Compute horizontal displacement penalty
1976     horizontal_displacement = torch.sum(torch.abs(torso_position[:,
1977     :2] - prev_torso_position[:, :2]), dim=-1)
1978     horizontal_displacement_penalty = torch.exp(-
1979     horizontal_displacement / 0.1) # Penalize large movements with
1980     temperature
1981
1982     # Compute leg forces usage reward
1983     contact_force_usage = torch.sum(torch.abs(
1984     right_leg_contact_force) + torch.abs(left_leg_contact_force),
1985     dim=-1)
1986     contact_force_usage_penalty = torch.exp(-contact_force_usage /
1987     10.0) # Penalize high contact force usage with temperature
1988
1989     # Compute joint usage reward (encourages proper leg extension
1990     and contraction)
1991     leg_joints_indices = torch.tensor([6, 7, 8, 12, 13, 14], device
1992     =device) # Indices of leg joints
1993     leg_joint_usage = torch.mean(dof_pos[:, leg_joints_indices],
1994     dim=-1)
1995     leg_joint_usage_reward = torch.exp(-torch.abs(leg_joint_usage)
1996     / 0.1) # Encourage movements from neutral position
1997
1998     # Sum all rewards and penalties
1999     total_reward = vertical_movement_reward +
2000     horizontal_displacement_penalty + contact_force_usage_penalty +
2001     leg_joint_usage_reward
2002
2003     # Create a dictionary for individual reward components
  
```

1998
 1999
 2000
 2001
 2002
 2003
 2004
 2005
 2006
 2007
 2008
 2009
 2010
 2011
 2012
 2013
 2014
 2015
 2016
 2017
 2018
 2019
 2020
 2021
 2022
 2023
 2024
 2025
 2026
 2027
 2028
 2029
 2030
 2031
 2032
 2033
 2034
 2035
 2036
 2037
 2038
 2039
 2040
 2041
 2042
 2043
 2044
 2045
 2046
 2047
 2048
 2049
 2050
 2051

```

reward_components = {
    'vertical_movement_reward': vertical_movement_reward,
    'horizontal_displacement_penalty':
horizontal_displacement_penalty,
    'contact_force_usage_penalty': contact_force_usage_penalty,
    'leg_joint_usage_reward': leg_joint_usage_reward
}

return total_reward, reward_components

```

HumanoidJump Task: Reward Function of Iteration 2

```

def compute_reward(
    torso_position: torch.Tensor,
    prev_torso_position: torch.Tensor,
    velocity_world: torch.Tensor,
    right_leg_contact_force: torch.Tensor,
    left_leg_contact_force: torch.Tensor,
    dof_pos: torch.Tensor
) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
    # Ensure all tensors are on the same device
    device = torso_position.device

    # Compute vertical torso movement reward
    vertical_movement = torso_position[:, 2] - prev_torso_position
   [:, 2]
    vertical_movement_reward = torch.clamp(vertical_movement, min
=0.0) # Reward positive vertical movement
    temperature_vertical = 0.05
    vertical_movement_reward = torch.exp(vertical_movement_reward /
temperature_vertical) # Use exponential scaling

    # Compute horizontal displacement penalty
    horizontal_displacement = torch.sum(torch.abs(torso_position[:,
:2] - prev_torso_position[:, :2]), dim=-1)
    temperature_horizontal = 0.1
    horizontal_displacement_penalty = torch.exp(-
horizontal_displacement / temperature_horizontal) # Penalize
large movements with temperature

    # Compute leg forces usage reward
    contact_force_usage = torch.sum(torch.abs(
right_leg_contact_force) + torch.abs(left_leg_contact_force),
dim=-1)
    temperature_force = 1.0
    contact_force_usage_penalty = torch.exp(-contact_force_usage /
temperature_force) # Penalize high contact force usage with
temperature

    # Compute joint usage reward (encourages proper leg extension
and contraction)
    leg_joints_indices = torch.tensor([6, 7, 8, 12, 13, 14], device
=device) # Indices of leg joints
    leg_joint_usage = torch.mean(dof_pos[:, leg_joints_indices],
dim=-1)
    temperature_joints = 0.05
    leg_joint_usage_reward = torch.exp(-torch.abs(leg_joint_usage)
/ temperature_joints) # Encourage movements from the neutral
position

```

```

2052
2053 # Sum all rewards and penalties
2054 total_reward = vertical_movement_reward +
2055 horizontal_displacement_penalty + contact_force_usage_penalty +
2056 leg_joint_usage_reward
2057
2058 # Create a dictionary for individual reward components
2059 reward_components = {
2060     'vertical_movement_reward': vertical_movement_reward,
2061     'horizontal_displacement_penalty':
2062     horizontal_displacement_penalty,
2063     'contact_force_usage_penalty': contact_force_usage_penalty,
2064     'leg_joint_usage_reward': leg_joint_usage_reward
2065 }
2066
2067 return total_reward, reward_components

```

HumanoidJump Task: Reward Function of Iteration 3

```

2070 def compute_reward(
2071     torso_position: torch.Tensor,
2072     prev_torso_position: torch.Tensor,
2073     velocity_world: torch.Tensor,
2074     right_leg_contact_force: torch.Tensor,
2075     left_leg_contact_force: torch.Tensor,
2076     dof_pos: torch.Tensor
2077 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
2078     # Ensure all tensors are on the same device
2079     device = torso_position.device
2080
2081     # Compute vertical torso movement reward
2082     vertical_movement = torso_position[:, 2] - prev_torso_position
2083    [:, 2]
2084     vertical_movement_reward = torch.clamp(vertical_movement, min
2085     =0.0) # Reward positive vertical movement
2086     temperature_vertical = 0.05
2087     vertical_movement_reward = torch.exp(vertical_movement_reward /
2088     temperature_vertical) # Use exponential scaling
2089
2090     # Compute horizontal displacement penalty
2091     horizontal_displacement = torch.sum(torch.abs(torso_position[:,
2092     :2] - prev_torso_position[:, :2]), dim=-1)
2093     temperature_horizontal = 0.05 # Adjusted temperature for finer
2094     tuning
2095     horizontal_displacement_penalty = torch.exp(-
2096     horizontal_displacement / temperature_horizontal) # Penalize
2097     large movements
2098
2099     # Compute leg forces usage reward
2100     contact_force_usage = torch.sum(torch.abs(
2101     right_leg_contact_force) + torch.abs(left_leg_contact_force),
2102     dim=-1)
2103     temperature_force = 5.0 # Adjusted to make contact force usage
2104     more noticeable
2105     contact_force_usage_penalty = torch.exp(-contact_force_usage /
2106     temperature_force) # Penalize high contact force usage
2107
2108     # Compute joint usage reward (encourages proper leg extension
2109     and contraction)
2110     leg_joints_indices = torch.tensor([6, 7, 8, 12, 13, 14], device
2111     =device) # Indices of leg joints

```

```

2106
2107     leg_joint_usage = torch.mean(dof_pos[:, leg_joints_indices],
2108                                 dim=-1)
2109     temperature_joints = 0.05
2110     leg_joint_usage_reward = torch.exp(-torch.abs(leg_joint_usage)
2111                                     / temperature_joints) # Encourage movements from the neutral
2112                                                         position
2113
2114     # Sum all rewards and penalties
2115     total_reward = vertical_movement_reward +
2116                 horizontal_displacement_penalty + contact_force_usage_penalty +
2117                 leg_joint_usage_reward
2118
2119     # Create a dictionary for individual reward components
2120     reward_components = {
2121         'vertical_movement_reward': vertical_movement_reward,
2122         'horizontal_displacement_penalty':
2123         horizontal_displacement_penalty,
2124         'contact_force_usage_penalty': contact_force_usage_penalty,
2125         'leg_joint_usage_reward': leg_joint_usage_reward
2126     }
2127
2128     return total_reward, reward_components

```

HumanoidJump Task: Reward Function of Iteration 4

```

2129
2130
2131 def compute_reward(
2132     torso_position: torch.Tensor,
2133     prev_torso_position: torch.Tensor,
2134     velocity_world: torch.Tensor,
2135     right_leg_contact_force: torch.Tensor,
2136     left_leg_contact_force: torch.Tensor,
2137     dof_pos: torch.Tensor
2138 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
2139     # Ensure all tensors are on the same device
2140     device = torso_position.device
2141
2142     # Compute vertical torso movement reward
2143     vertical_movement = torso_position[:, 2] - prev_torso_position
2144     vertical_movement_reward = torch.clamp(vertical_movement, min
2145     =0.0) # Reward positive vertical movement
2146     temperature_vertical = 0.04 # Adjusted temperature for better
2147     sensitivity
2148     vertical_movement_reward = torch.exp(vertical_movement_reward /
2149     temperature_vertical) # Use exponential scaling
2150
2151     # Compute horizontal displacement penalty
2152     horizontal_displacement = torch.sum(torch.abs(torso_position[:,
2153     :2] - prev_torso_position[:, :2]), dim=-1)
2154     temperature_horizontal = 0.1 # Increased temperature to
2155     penalize horizontal movement more
2156     horizontal_displacement_penalty = torch.exp(-
2157     horizontal_displacement / temperature_horizontal) # Penalize
2158     large movements
2159
2160     # Compute leg forces usage reward
2161     contact_force_usage = torch.sum(torch.abs(
2162     right_leg_contact_force) + torch.abs(left_leg_contact_force),
2163     dim=-1)

```

```

2160
2161 temperature_force = 0.1 # Significantly increase sensitivity
2162 to contact forces
2163 contact_force_usage_penalty = torch.exp(-contact_force_usage /
2164 temperature_force) # Penalize high contact force usage
2165
2166 # Compute joint usage reward (encourages proper leg extension
2167 and contraction)
2168 leg_joints_indices = torch.tensor([6, 7, 8, 12, 13, 14], device
2169 =device) # Indices of leg joints
2170 leg_joint_usage = torch.mean(dof_pos[:, leg_joints_indices],
2171 dim=-1)
2172 temperature_joints = 0.02 # Adjusted for joint usage
2173 sensitivity
2174 leg_joint_usage_reward = torch.exp(-torch.abs(leg_joint_usage)
2175 / temperature_joints) # Encourage movements from the neutral
2176 position
2177
2178 # Sum all rewards and penalties
2179 total_reward = vertical_movement_reward +
2180 horizontal_displacement_penalty + contact_force_usage_penalty +
2181 leg_joint_usage_reward
2182
2183 # Create a dictionary for individual reward components
2184 reward_components = {
2185     'vertical_movement_reward': vertical_movement_reward,
2186     'horizontal_displacement_penalty':
2187 horizontal_displacement_penalty,
2188     'contact_force_usage_penalty': contact_force_usage_penalty,
2189     'leg_joint_usage_reward': leg_joint_usage_reward
2190 }
2191
2192 return total_reward, reward_components

```

Humanoid Task: Reward Function of Iteration 5

```

2193
2194 def compute_reward(
2195     torso_position: torch.Tensor,
2196     prev_torso_position: torch.Tensor,
2197     velocity_world: torch.Tensor,
2198     right_leg_contact_force: torch.Tensor,
2199     left_leg_contact_force: torch.Tensor,
2200     dof_pos: torch.Tensor
2201 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
2202     # Ensure all tensors are on the same device
2203     device = torso_position.device
2204
2205     # Compute vertical torso movement reward
2206     vertical_movement = torso_position[:, 2] - prev_torso_position
2207    [:, 2]
2208     vertical_movement_reward = torch.clamp(vertical_movement, min
2209     =0.0) # Reward positive vertical movement
2210     temperature_vertical = 0.04 # Adjusted temperature for better
2211     sensitivity
2212     vertical_movement_reward = torch.exp(vertical_movement_reward /
2213     temperature_vertical) # Use exponential scaling
2214
2215     # Compute horizontal displacement penalty
2216     horizontal_displacement = torch.sum(torch.abs(torso_position[:,
2217     :2] - prev_torso_position[:, :2]), dim=-1)

```

```

2214
2215     temperature_horizontal = 0.05 # Decreased temperature for more
2216     sensitivity
2217     horizontal_displacement_penalty = torch.exp(-
2218     horizontal_displacement / temperature_horizontal) # Penalize
2219     large movements
2220
2221     # Compute leg forces usage penalty (Rewritten to reduce contact
2222     force)
2223     contact_force_usage = torch.sum(torch.abs(
2224     right_leg_contact_force) + torch.abs(left_leg_contact_force),
2225     dim=-1)
2226     temperature_force = 0.5 # Adjusted to penalize contact force
2227     usage
2228     contact_force_usage_penalty = torch.exp(-contact_force_usage /
2229     temperature_force) # Penalize high contact force usage
2230
2231     # Compute joint usage reward (encourages proper leg extension
2232     and contraction)
2233     leg_joints_indices = torch.tensor([6, 7, 8, 12, 13, 14], device
2234     =device) # Indices of leg joints
2235     leg_joint_usage = torch.mean(torch.abs(dof_pos[:,
2236     leg_joints_indices]), dim=-1)
2237     temperature_joints = 0.02 # Adjusted for joint usage
2238     sensitivity
2239     leg_joint_usage_reward = torch.exp(-leg_joint_usage /
2240     temperature_joints) # Encourage movements from the neutral
2241     position
2242
2243     # Sum all rewards and penalties
2244     total_reward = vertical_movement_reward +
2245     horizontal_displacement_penalty + contact_force_usage_penalty +
2246     leg_joint_usage_reward
2247
2248     # Create a dictionary for individual reward components
2249     reward_components = {
2250     'vertical_movement_reward': vertical_movement_reward,
2251     'horizontal_displacement_penalty':
2252     horizontal_displacement_penalty,
2253     'contact_force_usage_penalty': contact_force_usage_penalty,
2254     'leg_joint_usage_reward': leg_joint_usage_reward
2255     }
2256
2257     return total_reward, reward_components

```

Humanoid Task: Reward Function of Iteration 6

```

2256 def compute_reward(
2257     torso_position: torch.Tensor,
2258     prev_torso_position: torch.Tensor,
2259     velocity_world: torch.Tensor,
2260     right_leg_contact_force: torch.Tensor,
2261     left_leg_contact_force: torch.Tensor,
2262     dof_pos: torch.Tensor
2263 ) -> Tuple[torch.Tensor, Dict[str, torch.Tensor]]:
2264     # Ensure all tensors are on the same device
2265     device = torso_position.device
2266
2267     # Compute vertical torso movement reward
2268     vertical_movement = torso_position[:, 2] - prev_torso_position
2269    [:, 2]

```

```

2268
2269 vertical_movement_reward = torch.clamp(vertical_movement, min
2270 =0.0) # Reward positive vertical movement
2271 temperature_vertical = 0.03 # Fine-tuned temperature for
2272 better sensitivity
2273 vertical_movement_reward = torch.exp(vertical_movement_reward /
2274 temperature_vertical) # Use exponential scaling
2275
2276 # Compute horizontal displacement penalty
2277 horizontal_displacement = torch.sum(torch.abs(torso_position[:,
2278 :2] - prev_torso_position[:, :2]), dim=-1)
2279 temperature_horizontal = 0.04 # Decreased temperature for more
2280 sensitivity
2281 horizontal_displacement_penalty = torch.exp(-
2282 horizontal_displacement / temperature_horizontal) # Penalize
2283 large movements
2284
2285 # Compute leg forces usage penalty (encourage minimal contact
2286 force)
2287 contact_force_usage = torch.sum(torch.abs(
2288 right_leg_contact_force) + torch.abs(left_leg_contact_force),
2289 dim=-1)
2290 temperature_force = 0.5 # Adjusted to penalize contact force
2291 usage
2292 contact_force_usage_penalty = torch.exp(-contact_force_usage /
2293 temperature_force) # Penalize high contact force usage
2294
2295 # Compute joint usage reward (encourages proper leg extension
2296 and contraction)
2297 leg_joints_indices = torch.tensor([6, 7, 8, 12, 13, 14], device
2298 =device) # Indices of leg joints
2299 leg_joint_usage = torch.mean(torch.abs(dof_pos[:,
2300 leg_joints_indices]), dim=-1)
2301 temperature_joints = 0.02 # Fine-tuned for joint usage
2302 sensitivity
2303 leg_joint_usage_reward = torch.exp(-torch.abs(leg_joint_usage)
2304 / temperature_joints) # Encourage movements from the neutral
2305 position
2306
2307 # Sum all rewards and penalties
2308 total_reward = vertical_movement_reward +
2309 horizontal_displacement_penalty + contact_force_usage_penalty +
2310 leg_joint_usage_reward
2311
2312 # Create a dictionary for individual reward components
2313 reward_components = {
2314     'vertical_movement_reward': vertical_movement_reward,
2315     'horizontal_displacement_penalty':
2316 horizontal_displacement_penalty,
2317     'contact_force_usage_penalty': contact_force_usage_penalty,
2318     'leg_joint_usage_reward': leg_joint_usage_reward
2319 }
2320
2321 return total_reward, reward_components

```