

Improving LLM Code Reasoning via Semantic Equivalence Self-Play with Formal Verification

Anonymous ACL submission

Abstract

We introduce a self-play framework for semantic equivalence in Haskell, utilizing formal verification to guide adversarial training between a generator and an evaluator. The framework leverages Liquid Haskell proofs for validating equivalence and execution-based counterexamples for inequivalence, organized via a difficulty-aware curriculum. To facilitate this, we release **OpInstruct-HSx**, a synthetic dataset of $\approx 28k$ validated Haskell programs. Empirical experiments show that our evaluator transfers effectively to downstream tasks, achieving up to 13.3pp accuracy gain on EquiBench and consistent gains on PySecDB. Ablation studies on the SEQ-SINQ regimes indicate that while inequivalence supervision provides data volume, equivalence proofs are uniquely responsible for the model’s reasoning capabilities. The entire training pipeline and dataset are publicly released on GitHub and Hugging Face respectively.

1 Introduction

The rise of large language models (LLMs) has reshaped how software can be generated and maintained. Despite models such as Codex and Qwen2.5-coder have demonstrated strong capabilities in producing functional code from natural language prompts (Murphy et al., 2024; Hui et al., 2024), their outputs often fail to preserve the intended program behavior beyond basic test coverage (Laneve et al., 2025; Nguyen et al., 2025; Wei et al., 2025). This gap raises a fundamental challenge: How can we design training that explicitly teaches models to reason about semantic equivalence between programs? Addressing this problem is critical not only for reliable code generation but also for downstream applications such as program optimization, automated refactoring, and vulnerability detection.

Current approaches largely rely on test suites, which are insufficient for capturing deep semantic

properties and edge cases. To bridge this gap, we ground our framework in Haskell for three key reasons. First, its pure functional nature and strong static typing eliminate hidden state and side effects (Thompson, 2011) (See Appendix A for illustration), making equivalence reasoning mathematically tractable (Launchbury, 1993; Sestoft, 1997). Second, the Liquid Haskell ecosystem enables the generation of machine-checkable proofs via refinement types, making it possible to certify semantic equivalence for a subset of Haskell programs (Liquid Haskell Tutorial, 2025). This offers a source of formal supervision that is unavailable in languages such as Python or Java, where formal verification is much harder. Finally, training on underrepresented functional paradigms pushes the model beyond standard object-oriented patterns (van Dam et al., 2024; Giagnorio et al., 2025), encouraging deeper abstraction capabilities.

We propose a self-play framework for semantic equivalence, in which two specialized agents interact: Alice, a generator that produces variants of reference programs; and Bob, an evaluator trained to decide whether two programs are equivalent. The self-play loop alternates between program generation, verification through proofs or counterexamples, difficulty scoring, and fine-tuning of both agents. By framing the problem as a game between generator and evaluator, the system encourages progressively harder examples and deeper reasoning about semantics.

This work investigates three core questions regarding the utility of functional programming for LLM alignment. First, we examine the dynamics of self-play, asking whether an adversarial loop in a functional language can induce a progressive curriculum that improves semantic reasoning. Second, we evaluate cross-domain and cross-language transferability, assessing whether semantic reasoning skills acquired in Haskell generalize to zero-shot synthesis and vulnerability detection in broader

coding benchmarks. Finally, we perform a controlled ablation to determine the relative contributions of supervision signals, distinguishing between the effects of formal equivalence proofs and execution-based counterexamples on evaluator robustness.

2 Related Work

Determining semantic equivalence is extremely challenging in general and current LLMs often fail to recognise semantic equivalence in code (Laneve et al., 2025; Nguyen et al., 2025; Wei et al., 2025). By Rice’s Theorem, determining if two arbitrary programs are semantically equivalent is generally undecidable (Rice, 1953). Traditional approaches rely on unit testing; however, even extended test suites, such as HumanEval+ and MBPP+ remain insufficient to guarantee correctness (Gren and Antinyan, 2017; Chioteli et al., 2021). Similarly, symbolic execution offers path-sensitive analysis but suffers from combinatorial state-space explosion as program complexity increases (Badihi et al., 2020).

To address these incompleteness issues, recent research has pivoted towards formal verification. In the LLM domain, frameworks like DeepSeek-Prover (Xin et al., 2024a,b; Ren et al., 2025) and Kimina-Prover (Wang et al., 2025) have demonstrated that fine-tuning models on self-generated proofs (e.g., in Lean) significantly enhances their verifiable reasoning capabilities. However, generating fully machine-checkable equivalence proofs for imperative languages like Python or C++ is beyond the reach of current tools except for trivial cases (Miceli-Barone et al., 2025).

Our approach bridges this gap by using Haskell and Liquid Haskell. Liquid Haskell embeds refinement types into the language (Vazou et al., 2014), allowing logical properties to be verified automatically via Satisfiability Modulo Theories (SMT) solvers (Diatchki, 2015; Jhala et al., 2020; Liquid Haskell Tutorial, 2025). This framework enables the construction of machine-checkable lemmas, using reflection and Proof by Logical Evaluation (PLE), to formally certify pointwise equality between candidate functions. This provides a deterministic, high-fidelity feedback signal for the self-play loop.

Self-play has historically driven breakthroughs in game-playing agents like AlphaZero (Silver et al., 2017) and OpenAI Five (OpenAI et al., 2019).

In the coding domain, recent frameworks such as Sol-Ver (Lin et al., 2025) and AutoIF (Dong et al., 2025) adapt this by using model-generated unit tests and execution feedback to filter synthetic data, achieving significant gains on benchmarks like MBPP and IFEval.

Our work is most directly built upon the adversarial framework proposed by Miceli-Barone et al. (2025), known as the Semantic Inequivalence Game (SINQ). In their setup, a generator ("Alice") creates program variants and diverging inputs (counterexamples), while an evaluator ("Bob") attempts to detect inequivalence without seeing the generator’s justification. This adversarial loop provides a scalable curriculum that improves semantic reasoning. We extend this paradigm by including Semantic Equivalence tasks (SEQ). While Miceli-Barone et al. (2025) focused exclusively on inequivalence via execution feedback, we integrate formal verification using Liquid Haskell. This allows us to train on positive instances of equivalence supported by reasoning traces.

3 Methodology

We introduce the core methodology of the self-play framework for improving code reasoning in LLMs through two complementary tasks: the SEQ and the SINQ. In the SEQ task, the generator model is asked to produce a functionally identical variant of a given Haskell program, along with a formal proof certifying its equivalence. In contrast, the SINQ task requires generating a function that diverges from the original program on at least one input. Figure 1 shows the complete self-play framework, with difficulty-based supervised fine-tuning that guides the adaptive training loop.

3.1 Framework Overview

The self-play framework is formulated as a two-agent adversarial game between Alice (the generator) and Bob (the evaluator).

1. Each round starts with a reference Haskell program P , upon which Alice’s task is to generate program Q , which is either a semantically equivalent variant to P (with a proof) or a semantically inequivalent program (with a diverging input that demonstrates their different behaviour).
2. The proof or the diverging input is verified afterwards, and if it fails, Alice loses.

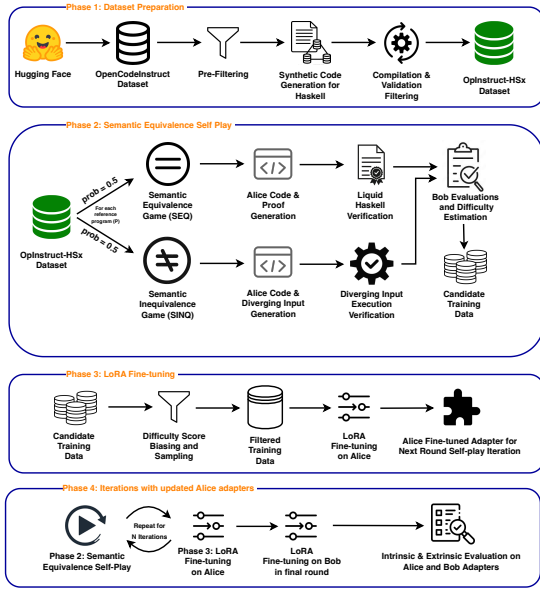


Figure 1: Overview of the semantic self-play framework for improving code reasoning in LLMs via Haskell.

- Bob then decides whether (P, Q) are semantically equivalent. If Bob’s judgment is accurate, he wins and Alice loses, vice versa.

Alice’s objective is to create instances that are difficult for Bob to classify, whereas Bob’s task is to correctly assess them. Through repeated interactions, both agents gradually improve their performance. Miceli-Barone et al. (2025) have proved that this adversarial framework has **no theoretical upper bound on model’s performance improvement**, and in principle both agents can learn endlessly about the complex programming logic while training on a real-world coding dataset.

3.2 Dataset Generation and Preparation

The availability of high-quality Haskell datasets remains extremely limited. Among the few usable resources, the most substantial one is the Blastwind dataset¹, which aggregates real-world source files scraped from public GitHub repositories. However, its utility is hindered by substantial noise: the data contains unannotated, inconsistently formatted, and often non-compilable code.

To address the scarcity of high-quality Haskell datasets, we adopt a complementary strategy: **introducing OpInstruct-HSx, where we generate a synthetic Haskell dataset** by adapting from the

¹<https://huggingface.co/datasets/blastwind/github-code-haskell-file>

nvidia-OpenCodeInstruct dataset (Ahmad et al., 2025), a large-scale, high-quality instruction corpus originally built for Python code generation.

The programs are first pre-filtered and then transformed into Haskell programs using the DeepSeek-R1-Distill-Llama-70B model. This process results in a synthetic dataset of Haskell programs. To ensure its quality, we applied an automated filtering and validation stage. For each generated Haskell program, we extract the function name and its argument types using syntactic heuristics, and synthesize a type-correct input using a recursive literal generator supporting common base types (e.g., Int, Bool, List, Tuple). Each program is compiled with Glasgow Haskell Compiler (GHC) and executed on the synthesised input. Only those that compile successfully and execute without errors are retained. This filtering process eliminates malformed or non-functional code, ensuring that the resulting dataset consists of minimally functional and executable Haskell programs.

Figure 2 shows the entire multi-stage filtering mechanism. We have contributed OpInstruct-HSx, a clean and executable Haskell dataset for both SEQ and SINQ games, which consists of approximately 28,000 validated Haskell functions derived from real-world problems. This dataset is made publicly available², serving as a high-quality synthetic Haskell resource for training LLMs in semantic reasoning tasks. The code to create the data and replicate the experiment is released on a public repository³.

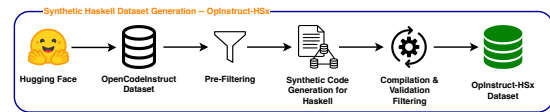


Figure 2: Full Pipeline for the OpInstruct-HSx dataset generation

3.3 The Self-Play Loop: Alice and Bob

Step 1: Program Selection and Branching

Let \mathcal{D} be the dataset of reference Haskell programs. We randomly choose a reference program $P \in \mathcal{D}$ and then select the SEQ game with 50% probability, otherwise we choose the SINQ game.

²URL redacted to preserve anonymity.

³URL redacted to preserve anonymity.

246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282

Step 2a: SEQ Game (Alice’s Turn)

In the SEQ game, Alice receives P and must synthesize Q such that for all inputs x :

$$\forall x \in \mathcal{X}. P(x) = Q(x)$$

To challenge Bob, Alice is encouraged to construct a hard instance Q , aiming for a maximum target difficulty level ($d=10$, defined in Section 3.4.1). In addition, Alice is required to produce a formal proof (in Liquid Haskell) of this semantic equivalence. See Appendix B for an SEQ instance.

Step 2b: SING Game (Alice’s Turn)

Alternatively in the SING game, Alice is instructed to produce a function Q that diverges from P on at least one input:

$$\exists x^* \in \mathcal{X} : P(x^*) \neq Q(x^*)$$

Alice is also incentivized to construct a difficult function Q that Bob is likely to misclassify, again targeting a maximum difficulty level ($d=10$). Alice must also output a diverging input x_a showing this inequivalence such that $P(x_a) \neq Q(x_a)$. See Appendix B for an SING instance.

Step 3: Verification through Liquid Haskell or Execution

- **SEQ Game:** Alice’s proof is verified by Liquid Haskell, which acts as an external oracle. If the proof is accepted, the proof is retained as a fine-tuning example.
- **SING Game:** The candidate x_a is tested by an execution: If $P(x_a) \neq Q(x_a)$, the instance is accepted. Otherwise, the sample is discarded.

All candidates undergo compilation, execution, and formal verification checks. This ensures the training data for both Alice and Bob remains high-quality and executable.

Step 4: Bob’s Turn – Difficulty Estimation

After Alice produces her candidate program Q , Bob is presented with both P and Q only and must decide whether the two programs are semantically equivalent. To estimate the challenge posed by each example, Bob is sampled N times and the proportion of correct responses from Bob is then used to compute a difficulty score. Further details on the difficulty-based curriculum and dataset sampling can be found in Section 3.4.

3.4 Implementation with Supervised Fine-Tuning with Difficulty Score

Given the practical challenges of reinforcement learning (Appendix C), we instead adopt rejection sampling supervised fine-tuning (SFT). We construct fine-tuning datasets by having Alice generate challenging programs and Bob learn from his own correct identifications. This semi-adversarial pipeline ensures Alice continually refines her ability to craft borderline-difficult SEQ or SING program pairs, while Bob steadily improves at its reasoning ability in semantic equivalence. Detailed prompt formats for both agents are provided in Appendix D.

3.4.1 Alice’s Training Data Selection

After Alice generates the candidate pairs (P, Q) , Alice’s outputs are not immediately used for fine-tuning. Instead, for each pair, Bob is asked to evaluate their semantic relation multiple times (typically $N = 10$). The number of correct Bob responses N_{success} determines the **difficulty score** \hat{d} as defined in Equation 1.

$$\hat{d} = d(P, Q) = 10 \times \left(1 - \frac{N_{\text{success}}}{N}\right) \quad (1)$$

However, most of Alice’s early generations are trivial for Bob. Including all examples would flood the dataset with low-difficulty cases that Bob already solves easily. Following the design in Miceli-Barone et al. (Miceli-Barone et al., 2025), we only retain examples that are sufficiently challenging, as determined by the difficulty score \hat{d} . Hence, shown in Schema 2, we split all (P, Q) pairs into $\mathcal{D}_{\text{Hard}}$ and $\mathcal{D}_{\text{Easy}}$, where τ is a chosen difficulty threshold (e.g., $\tau = 5$).

$$\begin{aligned} \mathcal{D}_{\text{Hard}} &= \{(P, Q) \mid d(P, Q) > \tau\}, \\ \mathcal{D}_{\text{Easy}} &= \{(P, Q) \mid d(P, Q) \leq \tau\} \end{aligned} \quad (2)$$

Finally, Alice’s training dataset is comprised of three SFT examples for each validated pairs (P, Q) . The first example pairs the prompts and Alice’s generation, and directly trains Alice to generate a challenging program Q , which helps improve Alice’s ability to craft precise SEQ or SING code (Schema 3). SP and UP denote the system and user prompts conditioned on the reference program P and a target difficulty $d = 10$. O represents the model’s raw output, which includes the chain-of-thought followed by the generated program Q (and

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
324
325
326
327

the diverging input x in the case of SING).

$$\begin{aligned} \text{SEQ} &= (\text{SP}_A^{eq}, \text{UP}_A^{eq}(P, 10), O_A^{eq}) \\ \text{SING} &= (\text{SP}_A^{inq}, \text{UP}_A^{inq}(P, 10), O_A^{inq}) \end{aligned} \quad (3)$$

We then select every hard example plus 20% as many easy ones sampled round-robin across integer difficulty bins to maintain a balanced curriculum.

The second example, designated as a ‘‘difficulty-prediction’’ instance, teaches Alice to self-assess the hardness of its own creations: by taking Alice’s generated program Q along with the Difficulty Prediction User Prompt and supervising on the numeric label ‘‘Difficulty level: \hat{d} ’’ (Schema 4). The training dataset is also biased towards hard examples, with easy examples comprising 50% of the subset. Alice learns to calibrate its difficulty estimates and ensures that its future generations are appropriately challenging.

$$\begin{aligned} \text{SEQ}_{\text{DP}} &= (\text{SP}_{A,DP}^{eq}, \text{UP}_{A,DP}^{eq}, \hat{d}) \\ \text{SING}_{\text{DP}} &= (\text{SP}_{A,DP}^{inq}, \text{UP}_{A,DP}^{inq}, \hat{d}) \end{aligned} \quad (4)$$

The third example concerns all (P, Q) pairs in the SEQ that have been successfully proved. We not only consider the generated candidate Q for a given P , but also the complete reasoning trace and valid Liquid Haskell proof script from Alice (Schema 5). These examples serve as supervised training signals to help Alice internalize the proof obligation π , corresponding to the refinement type $\forall x., P(x) = Q(x)$.

$$\text{SEQ}_{\text{proof}} = \left\{ (\text{SP}_A^{eq}, \text{UP}_A^{eq}(P), O_A^{eq}(P, Q, \pi)) \mid \text{LH}(P, Q) \vdash \pi \right\} \quad (5)$$

This three-part structure enables Alice to generate challenging variants, assess their difficulties, and internalize proof strategies.

3.4.2 Bob’s Training Data Selection

Bob’s training data is constructed from all of his correct (P, Q) pairs’ evaluations. Each training example (Schema 6) consists of the original pair, the system prompts, user prompts, and Bob’s response O_B . By sampling from a wide range of difficulties, Bob learns to reason and recognize both easy and subtle equivalence relations.

$$\mathcal{E}_{\text{Bob}} = (\text{SP}_B, \text{UP}_B(P, Q), O_B) \quad (6)$$

4 Experimental Setup

We use DeepSeek-R1-Distill-Qwen-7B as our base model for both Alice and Bob. Fine-tuning is performed using LoRA adaptors. Please see Appendix E for more details.

To investigate the distinct contributions of equivalence versus inequivalence supervision, we conduct the main experiment (E_0) alongside three controlled ablations (E_1 – E_3). These regimes systematically vary the game type probability and reference program budget (P) to isolate the impact of SEQ supervision under controlled data constraints. The configurations are shown in Table 1.

Regime	SEQ/SING	P	Goal
E_0 (Base)	50/50	500	Main
E_1	0/100	500	Max Volume
E_2	96/4	500	Balanced Yield
E_3	0/100	40	Vol. Control ($E_3 \approx E_2$)

Table 1: Experimental regimes to test SEQ impact.

The full codebase is available at the GitHub repository⁴ to generate the dataset OpInstruct-HSx, reproduce the experiments and evaluate model performance.

5 Results

The following evaluation results are structured to address the three research questions outlined in the introduction.

5.1 Intrinsic Evaluation

The goal of E_0 is teaching Alice (the generator) to generate increasingly challenging instances while Bob (the evaluator) to improve at judging program semantics.

5.1.1 Capability Assessment

As shown in Figure 3, the mean difficulty rises from 0.50 to 1.18 by round 7, indicating that Alice is crafting harder instances for a constant Bob. See Appendix H.1.1 for the difficulty trajectories breakdown for both SEQ and SING games.

In addition, we measure how much the evaluator model (Bob) improves after its first and only training round in the semantic equivalence self-play. Challenge instances are generated by the final trained generator model, Alice from round 7, using source programs from the training split of the OpInstruct-HSx and from the unseen test split. The

⁴URL redacted to preserve anonymity.

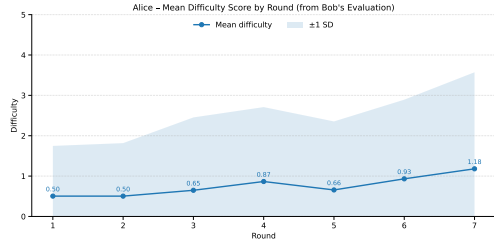


Figure 3: Mean and standard deviation for the difficulty scores for Alice’s generated instances from a fixed untrained Bob over 7 rounds.

resulting accuracies, summarised in Table 2, show modest improvements on unseen data.

Benchmark	Accuracy (%)		
	Base	Trained	Δ
OpInstruct-HSx (Train)	87.57	91.34	+3.77
OpInstruct-HSx (Test)	88.24	88.79	+0.55

Table 2: OpInstruct-HSx accuracy results comparing the Base Model and Trained (Bob adapter) models.

5.2 Extrinsic Evaluation

Having validated the agents’ improvement within the self-play loop, we now investigate whether this training enhances the model’s understanding of program semantics across different programming languages and domains.

5.2.1 Haskell Program Generation

In MBPP and HumanEval from MultiPL-E (Casano et al., 2022), both models’ compilation errors decrease substantially and Pass@1 improves on the Haskell tasks (shown in Table 3). Performance trajectories for Alice across self-play rounds are provided in Appendix H.1.2 and H.1.3.

The results indicate that semantic reasoning training transfers positively to Haskell code generation performance and robustness for both agents. For Bob, the gains confirm that the semantic discrimination skills learned during self-play can yield broader benefits in practical programming scenarios. And for Alice, the improvements suggest that training on high-quality, verifiable program transformations in self-play can enhance downstream synthesis accuracy and reduce compilation errors.

5.3 Coding Related Tasks Evaluation

We would like to investigate whether richer type semantics learnt in the evaluator model Bob can further enhance his ability to analyse correct and

semantically meaningful code across other coding paradigms, not just the functional ones.

5.3.1 EquiBench

The EquiBench (Wei et al., 2025) benchmark consists of six data categories: DCE (C pairs with dead/live code variations), x86-64 (identical assembly sequences via superoptimization), and CUDA (equivalent kernels with different scheduling) form the low-level systems group. The algorithmic group uses Python pairs from competitive programming, featuring OJ_A (algorithmic refactors), OJ_V (variable renaming), and OJ_VA (combined algorithmic and variable changes).

In the zero-shot prompting evaluation, our fine-tuned Bob model shows clear gains over the base model in several dataset configurations.

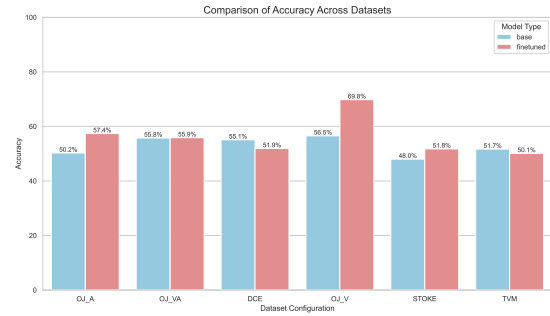


Figure 4: Accuracy comparison between base and fine-tuned Bob models across EquiBench dataset configurations.



Figure 5: F1 score comparison between base and fine-tuned Bob models across EquiBench dataset configurations.

From Figure 4 and Figure 5, the model transfers best to OJ_A and OJ_V because both categories (algorithmic refactors and variable renaming respectively) align with Bob’s learned skill on detecting high-level behavioural divergence, resulting in large gains in performance. OJ_VA’s performance stays modest as the authors stated that it is harder

Benchmark	Agent	Pass@1 (%)			Compilation Errors		
		Base	Trained	Δ	Base	Trained	Δ
HumanEval	Bob	17.7	26.4	+8.7	130	110	-20
HumanEval	Alice	17.7	26.3	+8.6	130	104	-26
MBPP	Bob	26.7	36.9	+10.2	250	203	-47
MBPP	Alice	26.7	34.3	+7.6	250	218	-32

Table 3: HumanEval and MBPP results (zero-shot prompting in Haskell), comparing the Base Model (untrained) and Trained models. Averages over 16 trials.

458 than pure renaming and closer to the “non-local
459 structural” end.

460 In contrast, DCE, STOKE, and TVM require
461 reasoning about low-level or non-functional seman-
462 tics absent from our Haskell curriculum. This mis-
463 match in language, abstraction level, and equiva-
464 lence definition explains the weaker or negligible
465 improvements on these tasks.

466 5.3.2 PySecDB

467 PySecDB is the first comprehensive dataset of
468 security-related commits in Python (Sun et al.,
469 2023). In Table 4, our fine-tuned Bob model
470 achieves consistent improvements over the base
471 model across all evaluated metrics. The results
472 indicate that Bob’s Haskell semantic reasoning
473 training transfers positively to Python vulnerability
474 detection, improving in identification of security-
475 relevant code changes.

Metric	Scores (%)		
	Base Model	Trained	Δ
Accuracy	67.6	68.8	+1.2
Precision	48.0	49.7	+1.7
Recall	55.1	59.2	+4.1
F1 Score	51.3	54.0	+2.7

Table 4: PySecDB vulnerability detection results comparing the Base Model and Trained (Bob adapter).

476 5.3.3 CodeXGlue

477 CodeXGLUE’s defect detection dataset (Zhou
478 et al., 2019) is a curated collection of over 27,000
479 C-language functions, each manually labeled to
480 indicate whether it contains a security-relevant de-
481 fect.

482 In Table 5, both the base and fine-tuned Bob
483 models perform near the 50% mark for accuracy,
484 which is indistinguishable from a random model.
485 CodeXGLUE’s dataset gives Bob only one C func-
486 tion and asks if it’s vulnerable. The vulnerabilities

Metric	Scores (%)		
	Base Model	Trained	Δ
Accuracy	48.8	50.5	+1.7
Precision	45.0	45.7	+0.7
Recall	51.1	41.7	-9.4
F1 Score	47.9	43.2	-4.7

Table 5: CodeXGLUE defect detection results comparing the Base Model and Trained (Bob adapter).

487 are low-level memory issues that require track-
488 ing pointers, memory allocation, and data lay-
489 outs. However, Haskell is memory-safe, garbage-
490 collected, and doesn’t have manual frees or pointer
491 arithmetic. Therefore, the trained model shows no
492 meaningful advantage over the base model on this
493 task.

494 5.4 Controlled Comparison of SEQ–SINQ 495 Regimes

496 A critical finding from E_0 is the asymmetry in num-
497 ber of validated programs between SINQ and SEQ
498 even under a uniform sampling policy in choos-
499 ing the game. In Figure 6, we have hundreds of
500 accepted SINQ instances per round, but very few
501 SEQ instances (single digits in most rounds). A
502 review of the output from the self-play shows that
503 while Alice can readily produce SINQ pairs, the
504 small reasoning model often struggles to produce
505 Liquid Haskell proofs for SEQ, which drastically
506 limits the number of compiled and verified SEQ
507 examples available for training. This imbalance
508 limits the diversity of program types in the training
509 buffer, with the current SEQ+SINQ configuration
510 heavily skewed with SINQ examples.

511 This subsection investigates whether incorpo-
512 rating SEQ supervision yields unique reasoning
513 benefits beyond those of SINQ alone, aiming to
514 isolate the impact of the supervision type from the
515 disparity in verified data volume (see Appendix F
516 for derivation details).

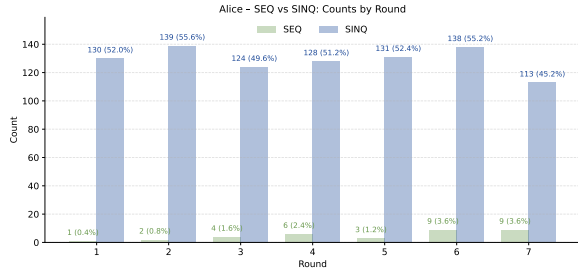


Figure 6: Counts and proportions of validated Alice generations by round.

5.4.1 Evaluation Results

We analyze the four experimental regimes to decouple the effects of supervision type (SEQ vs. SINQ) from training data volume. All results are presented in Appendix G and H.

Value of SEQ Supervision

Despite having significantly fewer verified training pairs in E_0 than in the pure SINQ regime E_1 (937 vs. 1,839), E_0 achieves superior performance on semantic equivalence tasks. This indicates that SEQ supervision improves high-level semantic judgment even when data volume is reduced. Furthermore, when strictly controlling for data volume (E_2 vs. E_3 , both ≈ 150 pairs), the inclusion of SEQ (E_2) yields consistent advantages on structural reasoning tasks over pure SINQ (E_3), shown in Table 7. This confirms that SEQ supervision confers a unique benefit that cannot be replicated by inequivalence signals alone.

Volume Trade-offs

While SEQ is qualitatively valuable, data volume remains critical. Regime E_2 , which aggressively prioritizes SEQ (96% attempts) to force a balanced validated distribution, suffers a sharp drop in total verified pairs ($N = 140$) due to low proof yields, causing performance degradation across all benchmarks compared to E_0 . The 50/50 mixed regime (E_0) therefore represents an acceptable trade-off, balancing the semantic depth of SEQ with the high verification yield of SINQ.

6 Conclusions

This paper investigates the use of self-play to improve code generation and semantic reasoning in LLMs, focusing on program equivalence in Haskell and Liquid Haskell. The framework introduces two agents: a generator of semantically equivalent or inequivalent program variants, and an evaluator

trained to judge equivalence. Evaluation shows that Bob benefits substantially: Haskell coding performance improves, and transfer is strong on high-level semantic reasoning tasks, but negligible on low-level or memory-oriented tasks. Controlled comparisons reveal that although inequivalence data yields far more verified pairs, the inclusion of equivalence supervision confers unique gains in semantic reasoning that cannot be achieved through inequivalence alone.

7 Future Work

There are several promising avenues for future research. First, extending training across more rounds could allow Bob to accumulate richer experience. Additionally, (1) scaling to larger parameter models and (2) employing full-parameter fine-tuning instead of LoRA may increase both Alice’s proof synthesis capabilities and Bob’s reasoning capacity.

A second priority is to improve the statistical robustness of the controlled comparisons. The experiments in regimes E_2 and E_3 were limited by small sample sizes, which may not be representative of the true differences between SINQ- and SEQ-based supervision. Running larger-scale experiments would provide stronger evidence about the relative contributions of the two regimes.

Furthermore, methodologically speaking, the Haskell self-play framework could be expanded to cover low-level program behaviors by introducing tasks that explicitly model memory usage, side effects, and bit-level operations.

Moreover, converting the self-play pipeline into a dedicated Haskell Equivalence Evaluation test set would help to provide a more direct and sustainable benchmark for semantic reasoning, especially given the scarcity of high-quality Haskell resources currently available.

Finally, reinforcement learning could be re-examined as a means of enhancing the SEQ game’s contribution. By running Alice long enough to generate sufficiently challenging programs and providing Bob with richer feedback through reward-based updates, it may be possible to overcome the current proof bottleneck and more effectively integrate equivalence reasoning into the self-play loop.

8 Limitations

A central limitation of this work lies in the proof synthesis bottleneck. Alice rarely produces Liquid

Haskell proofs that successfully pass PLE, resulting in very low validation yields for SEQ cases. This imbalance causes the majority of verified training data SING-dominated, reducing the contribution of SEQ supervision to Bob’s learning. The issue is likely exacerbated by the relatively small model size used in this study, which likely restricts Alice’s ability to generate structurally complex proofs.

Another limitation arises from the inherent constraints of Liquid Haskell itself. The verification process depends heavily on reflection and proof by logical evaluation, which are effective for local reasoning but struggle with certain classes of programs. Non-terminating behaviors, partial functions, and large-scale algebraic rewrites are difficult to certify, and in some cases impossible, within this fragment of the logic. Moreover, not all Haskell programs are reflectable, further narrowing the scope of tasks where equivalence can be formally verified.

Finally, there is a cross-domain mismatch in generalization. While Bob transfers strongly to high-level semantic reasoning tasks, such as OJ_A and OJ_V in EquiBench and vulnerability detection in PySecDB, his performance is notably weaker on low-level or stateful semantics, including DCE, STOKE, TVM, and CodeXGlue. This gap highlights a limitation of the current framework in addressing equivalence reasoning that depends on memory, side effects, or bit-level operations.

9 Ethical considerations

This work adversarially trains models to introduce hard to detect semantic modifications to computer programs and then to detect them, which may train them to detect security-relevant vulnerabilities or introduce obfuscated backdoors. In principle, these capabilities have a potential for malicious use, however, they can be also used for pro-social use to improve the reliability and security of code, whether generated by humans or LLMs. We believe that the net societal impact of models better able to reason about the subtleties of program semantics to be positive, especially as these capabilities are disseminated through Open Source releases, as we do in this work.

References

Wasi Uddin Ahmad, Aleksander Ficek, Mehrzad Samadi, Jocelyn Huang, Vahid Noroozi, Somshubra

Majumdar, and Boris Ginsburg. 2025. [Opencodeinstruct: A large-scale instruction tuning dataset for code llms](#). *Preprint*, arXiv:2504.04030.

Sahar Badihi, Faridah Akinotcho, Yi Li, and Julia Rubin. 2020. Ardif: scaling program equivalence checking via iterative abstraction and refinement of common code. In *Proceedings of the 28th ACM joint meeting on European software engineering conference and symposium on the foundations of software engineering*, pages 13–24.

Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael Greenberg, and Abhinav Jangda. 2022. [Multipl-e: A scalable and extensible approach to benchmarking neural code generation](#). *Preprint*, arXiv:2208.08227.

Efstathia Chioteli, Ioannis Batas, and Diomidis Spinellis. 2021. [Does unit-tested code crash? a case study of eclipse](#). In *25th Pan-Hellenic Conference on Informatics*, PCI 2021, page 260–264. ACM.

Iavor S. Diatchki. 2015. [Improving haskell types with smt](#). *SIGPLAN Not.*, 50(12):1–10.

Guanting Dong, Keming Lu, Chengpeng Li, Tingyu Xia, Bowen Yu, Chang Zhou, and Jingren Zhou. 2025. [Self-play with execution feedback: Improving instruction-following capabilities of large language models](#). In *The Thirteenth International Conference on Learning Representations*.

Alessandro Giagnorio, Alberto Martin-Lopez, and Gabriele Bavota. 2025. [Enhancing code generation for low-resource languages: No silver bullet](#). *arXiv preprint arXiv:2501.19085*.

Lucas Gren and Vard Antinyan. 2017. [On the relation between unit testing and code quality](#). In *2017 43rd Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*, page 52–56. IEEE.

Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, and 1 others. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.

Ranjit Jhala, Eric Seidel, and Niki Vazou. 2020. Programming with refinement types: An introduction to liquid haskell. <https://ucsd-progsys.github.io/liquidhaskell-tutorial/book.pdf>. Version 13, July 20, 2020.

Cosimo Laneve, Alvise Spanò, Dalila Ressi, Sabina Rossi, and Michele Bugliesi. 2025. [Assessing code understanding in llms](#). *Preprint*, arXiv:2504.00065.

John Launchbury. 1993. [A natural semantics for lazy evaluation](#). In *POPL’93: Proceedings of the 20th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 144–154. ACM.

705	Zi Lin, Sheng Shen, Jingbo Shang, Jason Weston, and Yixin Nie. 2025. Learning to solve and verify: A self-play framework for code and test generation . <i>Preprint</i> , arXiv:2502.14948.	759
706		760
707		761
708		762
709	Liquid Haskell Tutorial. 2025. Liquid haskell tutorial: Introduction. https://liquid.kosmik.us/org/01-intro.html . Accessed: 2025-07-18.	763
710		764
711		
712	Yihao Liu, Shuocheng Li, Lang Cao, Yuhang Xie, Mengyu Zhou, Haoyu Dong, Xiaojun Ma, Shi Han, and Dongmei Zhang. 2025. Superrl: Reinforcement learning with supervision to boost language model reasoning . <i>Preprint</i> , arXiv:2506.01096.	765
713		766
714		767
715		768
716		769
717	Matéo Mahaut and Francesca Franzon. 2025. Repetitions are not all alike: distinct mechanisms sustain repetition in language models. <i>arXiv preprint arXiv:2504.01100</i> .	770
718		771
719		772
720		773
721	Antonio Valerio Miceli-Barone, Vaishak Belle, and Ali Payani. 2025. Program semantic inequivalence game with large language models. <i>arXiv preprint arXiv:2505.03818</i> .	774
722		775
723		776
724		777
725	William Murphy, Nikolaus Holzer, Feitong Qiao, Leyi Cui, Raven Rothkopf, Nathan Koenig, and Mark Santolucito. 2024. Combining llm code generation with formal specifications and reactive program synthesis. <i>arXiv preprint arXiv:2410.19736</i> .	778
726		779
727		780
728		781
729		782
730	Thu-Trang Nguyen, Thanh Trong Vu, Hieu Dinh Vo, and Son Nguyen. 2025. An empirical study on capability of large language models in understanding code semantics. <i>Information and Software Technology</i> , page 107780.	783
731		784
732		785
733		786
734		787
735	OpenAI, :, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique P. d. O. Pinto, Jonathan Raiman, and 8 others. 2019. Dota 2 with large scale deep reinforcement learning . <i>Preprint</i> , arXiv:1912.06680.	788
736		789
737		790
738		791
739		792
740		793
741		794
742		795
743	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . <i>Preprint</i> , arXiv:2203.02155.	796
744		797
745		798
746		799
747		800
748		801
749		802
750		803
751	Z. Z. Ren, Zhihong Shao, Junxiao Song, Huajian Xin, Haocheng Wang, Wanxia Zhao, Liyue Zhang, Zhe Fu, Qihao Zhu, Dejian Yang, Z. F. Wu, Zhibin Gou, Shirong Ma, Hongxuan Tang, Yuxuan Liu, Wenjun Gao, Daya Guo, and Chong Ruan. 2025. Deepseek-prover-v2: Advancing formal mathematical reasoning via reinforcement learning for subgoal decomposition . <i>Preprint</i> , arXiv:2504.21801.	804
752		805
753		806
754		807
755		808
756		809
757		810
758		811
	Henry Gordon Rice. 1953. Classes of recursively enumerable sets and their decision problems. <i>Transactions of the American Mathematical society</i> , 74(2):358–366.	812
		813
	Peter Sestoft. 1997. Deriving a lazy abstract machine . <i>Journal of Functional Programming</i> , 7(3):231–264.	
	David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. 2017. Mastering chess and shogi by self-play with a general reinforcement learning algorithm . <i>Preprint</i> , arXiv:1712.01815.	
	Shiyu Sun, Shu Wang, Xinda Wang, Yunlong Xing, Elisa Zhang, and Kun Sun. 2023. Exploring security commits in python . <i>Preprint</i> , arXiv:2307.11853.	
	Simon Thompson. 2011. <i>Haskell: the craft of functional programming</i> . Addison-Wesley.	
	Tim van Dam, Frank van der Heijden, Philippe de Bekker, Berend Nieuwschepen, Marc Otten, and Maliheh Izadi. 2024. Investigating the performance of language models for completing code in functional programming languages: a haskell case study . <i>Preprint</i> , arXiv:2403.15185.	
	Niki Vazou, Eric L. Seidel, Ranjit Jhala, Dimitrios Vytiniotis, and Simon Peyton-Jones. 2014. Refinement types for haskell . <i>SIGPLAN Not.</i> , 49(9):269–282.	
	Haiming Wang, Mert Unsal, Xiaohan Lin, Mantas Baksys, Junqi Liu, Marco Dos Santos, Flood Sung, Marina Vinyes, Zhenzhe Ying, Zekai Zhu, and 1 others. 2025. Kimina-prover preview: Towards large formal reasoning models with reinforcement learning . <i>arXiv preprint arXiv:2504.11354</i> .	
	Anjiang Wei, Jiannan Cao, Ran Li, Hongyu Chen, Yuhui Zhang, Ziheng Wang, Yuan Liu, Thiago SFX Teixeira, Diyi Yang, Ke Wang, and 1 others. 2025. Equibench: Benchmarking large language models’ understanding of program semantics via equivalence checking . <i>arXiv preprint arXiv:2502.12466</i> .	
	Huajian Xin, Daya Guo, Zhihong Shao, Zhizhou Ren, Qihao Zhu, Bo Liu, Chong Ruan, Wenda Li, and Xiaodan Liang. 2024a. Deepseek-prover: Advancing theorem proving in llms through large-scale synthetic data . <i>arXiv preprint arXiv:2405.14333</i> .	
	Huajian Xin, Z. Z. Ren, Junxiao Song, Zhihong Shao, Wanxia Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, Wenjun Gao, Qihao Zhu, Dejian Yang, Zhibin Gou, Z. F. Wu, Fuli Luo, and Chong Ruan. 2024b. Deepseek-prover-v1.5: Harnessing proof assistant feedback for reinforcement learning and monte-carlo tree search . <i>Preprint</i> , arXiv:2408.08152.	
	Yaqin Zhou, Shangqing Liu, Jingkai Siow, Xiaoning Du, and Yang Liu. 2019. Devign: Effective vulnerability identification by learning comprehensive program	

814 semantics via graph neural networks. In *Advances*
815 *in Neural Information Processing Systems*, pages
816 10197–10207.

A Illustrative Haskell Features for Semantic Reasoning and Verification in LLM Self-Play

[1] Pure Semantics

Every function is a pure mapping from inputs to outputs, there is no hidden state or mutation. This purity ensures that two functions are semantically equivalent if they produce the same outputs for all inputs, regardless of how those outputs are computed.

```

1 -- revRec :: [a] -> [a]
2 revRec :: [a] -> [a]
3 revRec [] = []
4 revRec (x:xs) = revRec xs ++ [x]
5 -- revFold :: [a] -> [a]
6 revFold :: [a] -> [a]
7 revFold = foldl (flip (:)) []
8

```

Because Haskell is referentially transparent, `revRec xs` can be substituted with `revFold xs` wherever it appears. This makes reasoning about equivalence both feasible and formalizable.

[2] Static Typing and GHC Compile

Haskell's static type system provides strong guarantees at compile time. Once a program is accepted by the compiler, most classes of errors such as type mismatches and null de-referencing are eliminated. This makes the type checker an effective pre-filter for program validity in the self-play loop.

```

1 -- add :: Int -> Int -> Int
2 add x y = x + y
3

```

When writing `add True False`, GHC will raise a compilation error before running the program as `add` expects two `Int`.

B Liquid Haskell SEQ / SINQ Instances

Listing 1: Example of a SEQ Instance

```

-- Reference Function P
double :: Int -> Int
double x = x + x

-- Alice's generated Q (Semantic Equivalent)
double_alt :: Int -> Int
double_alt x = 2 * x

-- Liquid Haskell Proof
{-@ lemma_double_equiv :: x:Int
    -> { double x == double_alt x } @-}
lemma_double_equiv :: Int -> Proof
lemma_double_equiv x
= double x
=== double_alt x
*** QED

```

Listing 2: Example of a SINQ Instance

```

-- Original function P
sign :: Int -> String
sign n
| n < 0 = "negative"
| n == 0 = "zero"
| otherwise = "positive"

-- Alice's generated Q (semantically inequivalent)
signIneq :: Int -> String
signIneq n
| n <= 0 = "non-positive"
| otherwise = "positive"

```

```

-- Diverging input
x_a = 0
-- P x_a = sign 0 = "zero"
-- Q x_a = signIneq 0 = "non-positive"

```

C Implementation with Reinforcement Learning

Reinforcement learning (RL) offers a conceptually natural way to optimize the game by directly rewarding or penalizing generated programs based on the game outcome. A straightforward RL setup would proceed as follows based on the game rules:

C.1 RL Formulation

We treat Alice and Bob as two competing agents in a zero-sum Markov game.

- For the SINQ branch: An executor first checks that Q compiles and that $P(x_a) \neq Q(x_a)$ on Alice proposed diverging input x_a .
- For the SEQ branch: the SMT-based checker will execute the Liquid Haskell proof script asserting $\forall x. P(x) = Q(x)$.

If the above verification fails, the episode terminates with Alice receiving a negative reward $r_A = r_{\text{fail}} < 0$. Otherwise Bob observes state (P, Q) (but not Alice's diverging input x_a) and identify whether (P, Q) are semantically equivalent.

If Bob is correct, then Bob earns $r_B = r_{\text{success}} > 0$ and Alice is penalized with $r_A = r_{\text{too_easy}} < 0$. If Bob is incorrect, he gets $r_B = r_{\text{fail}} < 0$ and Alice wins $r_A = r_{\text{win}} > 0$. Both agents update their policies to maximize expected cumulative reward over many self-play episodes.

C.2 Potential Benefits of RL

RL can be layered on top of the semantic equivalence self-play by treating Alice's generator as an RL agent whose policy is updated based on proof outcomes. Over many episodes, this encourages the policy to avoid classes of common semantic errors (e.g. off-by-one edge cases) and refines its internal value function to distinguish subtle equivalence-breaking patterns, resulting in an LLM that is both more precise and robust.

C.3 Practical Challenges

- **Sparse, High-Variance Rewards:** In our reasoning tasks, reward is only provided at the

898 end of a multi-step chain of thought. As a re-
 899 sult, the LLMs often struggle to find success-
 900 ful trajectories, making it difficult for LLMs
 901 to learn which intermediate steps contributed
 902 to success (Liu et al., 2025).

- 903 • **Credit Assignment:** Attributing success or
 904 failure back to specific reasoning tokens or
 905 code-editing steps is nontrivial, and naive re-
 906 ward assignment can lead to undesirable short-
 907 cuts, such as inserting semantically neutral
 908 modifications to exploit reward signals with-
 909 out genuine reasoning.
- 910 • **Computational Resources Constraints:** RL
 911 is more computationally expensive than su-
 912 pervised fine-tuning due to the repeated sam-
 913 pling, model rollouts, and reward computation
 914 across multiple steps (Ouyang et al., 2022).

915 D Prompt Template

916 D.1 Alice SEQ System Prompt

```
917 You are a helpful and expert Haskell
918 programmer, powered by Liquid Haskell.
919 Your job is to *transform* any given function
920 `P` into a new function `Q` that:
921
922 - Is syntactically correct Haskell.
923 - Is semantically equivalent:  $\forall x. P\ x == Q\ x$ .
924 - Uses a *different* implementation --
925 Implement a non-structural change rather
926 than merely swapping operator order.
927 - Uses a different function name (e.g. add a
928 trailing `_alt`).
```

929 The original program and your program will be
 930 used in a test to evaluate the skill of an
 931 expert Haskell programmer who will
 932 identify whether your program Q is
 933 semantic equivalent to the original
 934 program P, so make sure that the
 935 difference you introduce are not very easy
 936 to understand.

937 You will be given a difficulty level from 0
 938 (easiest) to 10 (hardest) to target. E.g.
 939 difficulty level 0 means that an expert
 940 computer scientist in the bottom decile or
 941 above should be able to find a diverging
 942 example, difficulty level 9 means that
 943 only an expert computer scientist in the
 944 top decile should be able to find a
 945 diverging example, and difficulty level 10
 946 means that only the top 0.01 or less of
 947 expert Haskell programmer should be able
 948 to find a diverging example.

949 Always think through your transformation steps
 950 in ` $\langle \text{think} \rangle \dots \langle /\text{think} \rangle$ `, then emit exactly:

```
951 Generated Program `Q` :
952 ```haskell
953 <your Q here>
954 ```
```

963 D.2 Alice SEQ User Prompt

```
964 Difficulty level: {difficulty_level}
965 Here is the original Haskell function `P` :
966 ```haskell
967 {program_p_completion}
968 ```
969
970 Its argument type is
971 ```haskell
972 t = {t}
973 ```
974
975 Your task: produce a new function `Q` that
976 satisfies the system prompt requirements.
977 - Make sure `Q` has a different name (e.g.
978 append a `_alt`).
```

979 - Avoid trivial symmetric rewrites - show a
 980 genuine alternative implementation.
 981 - Do not include any extra commentary beyond
 982 the required ` $\langle \text{think} \rangle \dots \langle /\text{think} \rangle$ ` and the
 983 `Generated Program `Q` :` block.
 984 - Where appropriate, feel free to use Prelude
 985 functions such as foldr, map, or zipWith
 986 to encourage diverse strategies.
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025
 1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079
 1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133
 1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187
 1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241
 1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295
 1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349
 1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403
 1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457
 1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619
 1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673
 1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727
 1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781
 1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835
 1836
 1837
 1838
 1839
 1840
 1841
 1842
 1843
 1844
 1845
 1846
 1847
 1848
 1849
 1850
 1851
 1852
 1853
 1854
 1855
 1856
 1857
 1858
 1859
 1860
 1861
 1862
 1863
 1864
 1865
 1866
 1867
 1868
 1869
 1870
 1871
 1872
 1873
 1874
 1875
 1876
 1877
 1878
 1879
 1880
 1881
 1882
 1883
 1884
 1885
 1886
 1887
 1888
 1889
 1890
 1891
 1892
 1893
 1894
 1895
 1896
 1897
 1898
 1899
 1900
 1901
 1902
 1903
 1904
 1905
 1906
 1907
 1908
 1909
 1910
 1911
 1912
 1913
 1914
 1915
 1916
 1917
 1918
 1919
 1920
 1921
 1922
 1923
 1924
 1925
 1926
 1927
 1928
 1929
 1930
 1931
 1932
 1933
 1934
 1935
 1936
 1937
 1938
 1939
 1940
 1941
 1942
 1943
 1944
 1945
 1946
 1947
 1948
 1949
 1950
 1951
 1952
 1953
 1954
 1955
 1956
 1957
 1958
 1959
 1960
 1961
 1962
 1963
 1964
 1965
 1966
 1967
 1968
 1969
 1970
 1971
 1972
 1973
 1974
 1975
 1976
 1977
 1978
 1979
 1980
 1981
 1982
 1983
 1984
 1985
 1986
 1987
 1988
 1989
 1990
 1991
 1992
 1993
 1994
 1995
 1996
 1997
 1998
 1999
 2000

993 D.3 Lemma SEQ Proof System Prompt

```
994 You are an expert Haskell/Liquid Haskell prover.
995 You are asked to prove that two reflected
996 functions are equivalent.
997
998 The most basic proof should be in the following
999 format:
1000 ```haskell
1001 {-@ lemma_{func_name_p}_equiv :: x:{arg_type}
1002   -> {{ {func_name_p} x == {func_name_q} x
1003     }} @-}}
1004 lemma_{func_name_p}_equiv :: {arg_type}
1005   -> Proof
1006   lemma_{func_name_p}_equiv x
1007     = {func_name_p} x
1008     == {func_name_q} x
1009     * QED
1010 ```
1011
1012 However, you should also use more advanced
1013 proof techniques if necessary.
1014
1015 Few-Shot Example 1:
1016
1017 ```haskell
1018 {-@ LIQUID "--reflection" @-}}
1019 {-@ LIQUID "--ple" @-}}
1020
1021 module MyTest where
1022
1023 import Language.Haskell.Liquid.ProofCombinators
1024
1025 -- Alice program P
1026 {-@ reflect double @-}}
1027 double :: Int -> Int
1028 double x = x + x
1029
1030 -- Alice proposes Q
1031 {-@ reflect double' @-}}
1032 double' :: Int -> Int
1033 double' x = 2 * x
1034
1035 -- Here is the full lemma, from annotation to
1036 QED:
1037 {-@ lemma_double_equiv :: x:Int -> {{ double x
1038   == double' x }} @- }}
1039 lemma_double_equiv :: Int -> Proof
1040 lemma_double_equiv x
```

```

1043 = double x
1044 === double' x
1045 * QED
1046 ```
1047
1048 Few-Shot Example 2:
1049 ```haskell
1050 {-@ LIQUID "--reflection" @-}
1051 {-@ LIQUID "--ple" @-}
1052 module Equiv where
1053
1054 import Language.Haskell.Liquid.ProofCombinators
1055
1056 {-@ reflect addNumbers @-}
1057 addNumbers :: Int -> Int -> Int
1058 addNumbers a b = a + b
1059
1060 {-@ reflect addNumbers' @-}
1061 addNumbers' :: Int -> (Int -> Int)
1062 addNumbers' a = \b -> a + b
1063
1064 -- Alice detailed proof of equivalence
1065 lemma_addNumbers_equiv :: Int -> Int -> Proof
1066 lemma_addNumbers_equiv x y
1067 = addNumbers x y
1068 === addNumbers' x y
1069 * QED
1070
1071 When you answer, output only the complete lemma
1072 block in the same style:
1073 1. Use the `{-@ lemma... @-}` annotation ,
1074 with the exact naming pattern
1075 lemma_<P>_equiv
1076 2. The Haskell type signature
1077 3. The function definition with `===` steps
1078 4. End with `* QED`
1079 5. Please put your proof between ```haskell and
1080 ```
1081
1082 No extra text, no additional comments.
1083 Your answer must match the example format
1084 exactly, without trailing whitespace or
1085 newlines outside the code block.

```

D.4 Lemma SEQ Proof User Prompt

```

1087 {error_msg_section}
1088 {equiv_code}
1089 -----
1090
1091 Your task: Produce the proof of equivalence for
1092 the following function:
1093 {func_name_p} x == {func_name_q} x` for all
1094 `x`.
1095
1096 ```haskell
1097 {-@ LIQUID "--reflection" @-}
1098 {-@ LIQUID "--ple" @-}
1099 module Equiv where
1100 import Language.Haskell.Liquid.ProofCombinators
1101
1102 {-@ reflect {func_name_p} @-}
1103 {program_p_content}
1104
1105 {-@ reflect {func_name_q} @-}
1106 {program_q_content}
1107
1108 -- Your complete proof of equivalence
1109 /* PROOF BODY HERE */
1110 ```
1111
1112 <think>

```

D.5 Alice SING System Prompt

```

1119
1120

```

```

You are an expert Haskell programmer. Your task
is to generate a semantically inequivalent
variant of a given Haskell program, which
means that there must exist at least a
diverging input example such that the
original program and your program either
produce different outputs or exceptions,
or one halts and the other one does not
halt.
You must also provide a diverging input, which
is a valid input for both programs, but on
which they produce different outputs.
A good inequivalent program `Q` should be
subtly different from `P`.
A good diverging input `x` should be simple and
clearly demonstrate the semantic
difference between `P` and `Q`.
The original program and your program will be
used in a test to evaluate the skill of an
expert Haskell programmer who will have to
produce a diverging example (not
necessarily the same as yours), so make
sure that the difference you introduce are
not very easy to understand.
You will be given a difficulty level from 0
(easiest) to 10 (hardest) to target. E.g.
difficulty level 0 means that an expert
computer scientist in the bottom decile or
above should be able to find a diverging
example, difficulty level 9 means that
only an expert computer scientist in the
top decile should be able to find a
diverging example, and difficulty level 10
means that only the top 0.01 or less of
expert Haskell programmer should be able
to find a diverging example.
First, think step-by-step and write down your
analysis of program `P` and your strategy
for creating an inequivalent program `Q`.
Enclose this reasoning within <think>
and </think> tags.
After the thinking block, the final answer
could only be in the following format,
without any additional explanation or
context.
Final output MUST be exactly:
Generated Program `Q`:
```haskell
<Your generated Haskell code for `Q`>
```
Diverging Input `x`:
<The diverging input `x`>

```

D.6 Alice SING User Prompt

```

Difficulty level: {difficulty_level}
Original program `P`:
```haskell
{program}
```
<think>

```

D.7 Bob System Prompt

```

You are an expert Haskell programmer. You are
given two Haskell programs, `P` and `Q`.
Your task is to determine if they are
semantically equivalent.

```

```

1202 Use the following format to respond:
1203 # Equivalent?
1204 Yes or No
1205
1206 If the programs are equivalent, respond with
1207 your thought process and a final output
1208 with:
1209 # Equivalent?
1210 Yes
1211
1212 If they are inequivalent, respond with your
1213 thought process and a final output with:
1214 # Equivalent?
1215 No

```

1218 D.8 Bob User Prompt

```

1220 Program `P`:
1221 ```haskell
1222 {program_p}
1223 ```
1224
1225
1226 Program `Q`:
1227 ```haskell
1228 {program_q}
1229 ```
1230 <think>
1231

```

1233 D.9 Alice SEQ Difficulty Prediction System Prompt

```

1237 Difficulty level: Any
1238 Program P
1239 ```haskell
1240 {program_p}
1241 ```
1242
1243 The program Q below is semantically equivalent
1244 to the original program P, where the
1245 equivalence is due to the fact that the
1246 two programs have the same behavior on all
1247 inputs.
1248 Program Q
1249 ```haskell
1250 {program_q}
1251 ```

```

1253 D.10 Alice SING Difficulty Prediction System Prompt

```

1257 Difficulty level: Any
1258 program P
1259 ```haskell
1260 {program_p}
1261 ```
1262
1263 The program Q below is semantically
1264 inequivalent to the original program P,
1265 where the inequivalence is due to the fact
1266 that the two programs have different
1267 behavior on some inputs.
1268 Program Q
1269 ```haskell
1270 {program_q}
1271 ```

```

1273 D.11 Alice Difficulty Prediction User Prompt

```

1274 Predict the difficulty level of the instance of
1275 Program Q compared to Program P. Just
1276 write "Difficulty level: D" where D is
1277 your prediction, do not write anything
1278 else.
1279
1280

```

1283 E Main Experiment Setup

1284 E.1 Model Configs and Experiment Settings

1285 We employ DeepSeek-R1-Distill-Qwen-7B as
1286 the base model. The decoding parameters are con-
1287 figured as: temperature $T = 0.6$, top- $p = 0.95$,
1288 top- $k = 20$, presence penalty = 1.5 (Applied to
1289 reduce repetition in outputs from smaller reason-
1290 ing models (Mahaut and Franzon, 2025)), and a
1291 context length of 32,768 tokens.

1292 In E_0 , Alice is configured with a 50/50 prob-
1293 ability of playing either the SEQ or SING game
1294 (Section 3.3), with the difficulty level set as 10 in
1295 her prompts (Appendix D.2 and D.6). Bob serves
1296 solely as an evaluator and difficulty scorer, and he
1297 is not fine-tuned iteratively but only once after all
1298 rounds are complete.

1299 E.2 Dataset Size and Difficulty Threshold

1300 Due to computational constraints, each experimen-
1301 tal run uses a maximum of 500 Haskell reference
1302 programs⁵, i.e., $|\mathcal{D}| = P = 500$. The difficulty
1303 threshold τ is set to 3 (rather than the default 5)
1304 to increase the number of training examples avail-
1305 able to Alice. For future experiments with greater
1306 resources, it is recommended to increase $|\mathcal{D}|$ and
1307 restore $\tau = 5$ to more closely simulate a fully ad-
1308 versarial setting.

1309 E.3 Fine-Tuning Protocol

1310 After each self-play round, Alice is fine-tuned on
1311 the combination of newly generated instances and
1312 the retained examples from all previous rounds.
1313 More importantly, each round’s fine-tuning is ini-
1314 tialized from the original base model rather than
1315 from the fine-tuned adapter checkpoint obtained
1316 in the previous round. This strategy is adopted to
1317 avoid bias accumulation across iterations and to
1318 maintain stable and unbiased difficulty estimation.

1319 During generation in round $i + 1$, Alice uses the
1320 LoRA learned in round i , but fine-tuning for that
1321 round still begins from the base model.

1322 E.4 Experimental Duration

1323 The main run comprises 7 self-play rounds. In each
1324 round, Alice is fine-tuned according to the above
1325 protocol. Bob undergoes a single fine-tuning step
1326 after all 7 rounds are complete. For both Alice and
1327 Bob, the fine-tuning minimises the loss only on the
1328 model’s own outputs.

⁵A main run with fine-tuning takes about 3 days on 4 NVIDIA L40S GPUs.

F Rationale and Derivations for the SEQ-SINQ Comparison Experiment

F.1 Fairness Criterion and Derivation

Let P be the number of reference programs attempted per round, $p \in [0, 1]$ the probability of playing SEQ (so $1 - p$ for SINQ), and let $r_{\text{SEQ}}, r_{\text{SINQ}}$ denote the verification yields (probability that an attempt becomes a verified training example). By linearity of expectation, the expected number of verified examples per round is

$$K(P, p) = P(p r_{\text{SEQ}} + (1 - p) r_{\text{SINQ}}).$$

To isolate the effect of SEQ vs. SINQ unbalanced volume size, we match the expected count K of verified pairs across arms.

F.2 Yield Estimation from the Main Run E_0

From the 7-round of main experiment E_0 , $|\mathcal{D}| = P = 500$ (per-round attempts ≈ 250 SEQ and 250 SINQ):

| Game | Total Validated | Attempts | Yield (\hat{r}) |
|------|-----------------|----------|---------------------|
| SEQ | 34 | 1,750 | 1.94% |
| SINQ | 903 | 1,750 | 51.6% |

Table 6: Validation yields from E_0 . Validated counts by round were: SEQ (1, 2, 4, 6, 3, 9, 9) and SINQ (130, 139, 124, 128, 131, 138, 113).

F.3 Predicting Supervision for Each Arm

To achieve a balanced distribution of verified training examples (approx. 50/50 split), we configured regime E_2 with a mix of 96% SEQ and 4% SINQ. With a reference budget of $P = 500$, the estimated verified yields per round are:

$$N_{\text{SEQ}} \approx 500 \cdot 0.96 \cdot 0.0194 \approx 9.3 \quad (7)$$

$$N_{\text{SINQ}} \approx 500 \cdot 0.04 \cdot 0.516 \approx 10.3 \quad (8)$$

Totaling these yields using $K(P, p)$ results in an expected $K \approx 19.6$ verified examples per round, achieving the desired parity between the two supervision signals.

F.4 Matched comparator for E_2

Choose P_{E_3} in a 100% SINQ run so that

$$\begin{aligned} P_{E_3} \cdot r_{\text{SINQ}} &= K_{E_2} \\ \Rightarrow P_{E_3} &= \frac{19.632}{0.516} = 38.04 \approx 40 \end{aligned} \quad (9)$$

F.5 What this controls, and what it does not

Matching K ensures each arm receives similar amounts of verified training, so outcome differences are attributable to the type of game (presence/absence of SEQ) rather than volume. However, residual differences may still arise from the difficulty distribution of accepted items and class-imbalance within Bob’s updates.

F.6 Comparison Objectives

The experimental design supports the following four key pairwise comparisons:

- E_0 vs. E_1 : Under fixed compute resources (constant P), does a mixed 50% SEQ / 50% SINQ regime yield greater performance improvements than a pure SINQ regime similar to Miceli et al. setting (Miceli-Barone et al., 2025), given that the latter produces substantially more verified training examples due to higher verification success rates in SINQ?
- E_0 vs. E_2 : With fixed compute resources, does shifting towards a more balanced number of verified SEQ and SINQ examples, at the cost of reducing the overall number of verified training pairs, lead to improved evaluator performance compared to a SINQ-dominated main experiment?
- E_1 vs. E_2 : Under the same compute constraints, does a pure SINQ regime (maximising verified training pairs) outperform a more balanced SEQ–SINQ dataset that sacrifices data volume for greater semantic diversity?
- E_2 vs. E_3 : When the number of verified training pairs is held constant, does the inclusion of SEQ supervision yield measurable performance benefits over an SINQ-only regime?

F.7 Summary of experiments

- E_0 (main): 50/50, $P = 500$.
- E_1 : 100% SINQ, $P = 500$.
- E_2 : 96% SEQ / 4% SINQ, $P = 500$.
- E_3 (matched to E_2): 100% SINQ, $P = 40$ to equalize expected verified pairs with E_2 .

This suite isolates the causal effect of adding SEQ supervision under controlled total verified pairs, enabling a fair assessment of whether SEQ contributes beyond SINQ alone.

| Benchmark | Metric | Score (%) | | | |
|-------------------|--------|----------------|----------------|----------------|----------------|
| | | E ₀ | E ₁ | E ₂ | E ₃ |
| HumanEval (Bob) | Pass@1 | 26.4 | 25.8 | 22.7 | 22.0 |
| HumanEval (Alice) | Pass@1 | 26.3 | 28.1 | 21.8 | 18.4 |
| MBPP (Bob) | Pass@1 | 36.9 | 38.8 | 34.7 | 32.4 |
| MBPP (Alice) | Pass@1 | 34.3 | 39.3 | 32.9 | 28.0 |
| EquiBench (OJ_A) | Acc | 57.4 | 54.0 | 55.5 | 53.3 |
| EquiBench (OJ_A) | F1 | 50.7 | 44.4 | 48.9 | 45.1 |
| EquiBench (OJ_V) | Acc | 69.8 | 61.4 | 63.2 | 65.8 |
| EquiBench (OJ_V) | F1 | 69.7 | 58.0 | 61.5 | 64.7 |
| EquiBench (OJ_VA) | Acc | 55.9 | 54.5 | 54.7 | 53.7 |
| EquiBench (OJ_VA) | F1 | 47.5 | 45.4 | 47.1 | 45.2 |
| EquiBench (STOKE) | Acc | 51.8 | 50.9 | 51.0 | 50.8 |
| EquiBench (STOKE) | F1 | 18.3 | 15.4 | 16.8 | 15.9 |
| EquiBench (TVM) | Acc | 50.1 | 50.3 | 48.8 | 50.8 |
| EquiBench (TVM) | F1 | 38.3 | 39.9 | 33.5 | 33.9 |
| EquiBench (DCE) | Acc | 51.9 | 52.6 | 51.5 | 51.3 |
| EquiBench (DCE) | F1 | 47.7 | 45.4 | 45.2 | 44.4 |
| PySecDB | Acc | 68.8 | 67.3 | 68.7 | 68.7 |
| PySecDB | F1 | 54.0 | 51.9 | 53.9 | 53.6 |
| CodeXGlue | Acc | 50.5 | 49.2 | 49.9 | 48.9 |
| CodeXGlue | F1 | 43.2 | 44.0 | 44.1 | 48.5 |

Table 7: Benchmark Performances over all experiments. Averages over 16 samples.

| Game | Count | | | |
|-------|----------------|----------------|----------------|----------------|
| | E ₀ | E ₁ | E ₂ | E ₃ |
| SEQ | 34 | – | 62 | – |
| SINQ | 903 | 1,839 | 78 | 156 |
| Total | 937 | 1,839 | 140 | 156 |

Table 8: Total verified generation from Alice of all experiments over 7 rounds.

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

H Extended Regime Results

H.1 Experiment E_0

Most of the evaluation results are included in the main report.

H.1.1 Difficulty trajectories for both task types

When we inspect the difficulty trajectories for both task types in Figure 7, they move upward across rounds.

- For SEQ, although the sample size is not statistically significant, its mean difficulty score jumps from near 0.0 in early rounds to 2.0 – 3.0 in later ones, signaling that Alice manages to construct are non-trivial from Bob’s perspective.
- For SING, its averages being greater than the medians (0.0 for all rounds) indicates a right-skewed distribution. This suggests that while most generated instances are of lower difficulty, Alice occasionally produces much harder examples that significantly influence the mean. Starting in round 4, the magnitude of the difference between the mean and median becomes larger, and by round 7 the mean even exceeds the upper quartile. Therefore, Alice is producing increasingly harder examples over successive rounds.

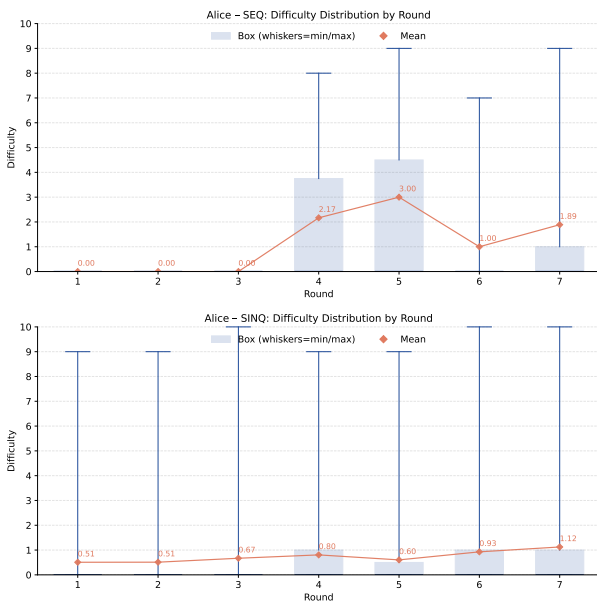


Figure 7: Mean and box-whiskey plots of the difficulty scores of Alice’s generated instances by round. Top: SEQ Game. Bottom SING Game.

Both SEQ and SING setups show rising difficulty scores across rounds, indicating that Alice improves in generating more challenging instances for Bob.

1427

1428

1429

1430

H.1.2 Alice’s HumanEval (Haskell) Performance

1431

1432

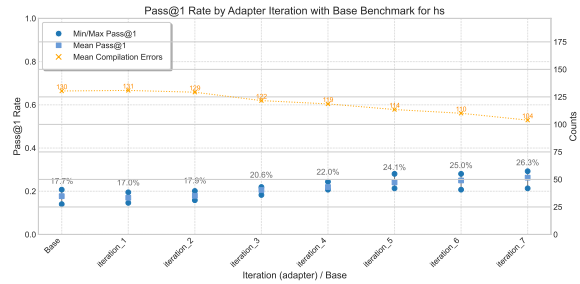


Figure 8: Alice’s HumanEval results. Averages over 16 trials.

H.1.3 Alice’s MBPP (Haskell) Performance

1433

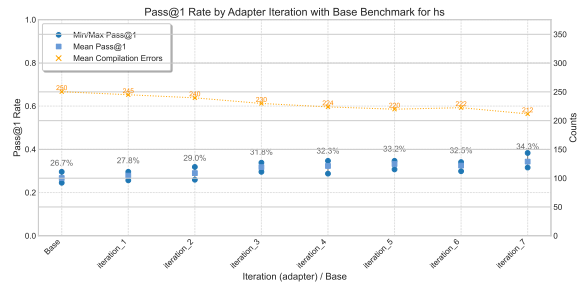


Figure 9: Alice’s MBPP results. Averages over 16 trials.

H.2 Experiment E_1

1434

H.2.1 Alice’s HumanEval (Haskell) Performance

1435

1436

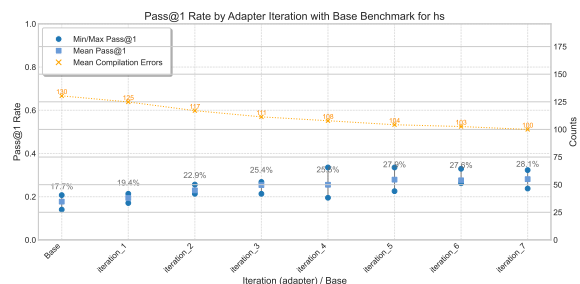


Figure 10: Alice’s HumanEval results. Averages over 16 trials.

H.2.2 Alice's MBPP (Haskell) Performance

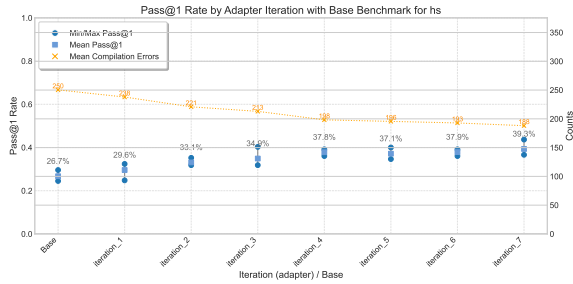


Figure 11: Alice's MBPP results. Averages over 16 trials.

H.3.2 Alice's MBPP (Haskell) Performance

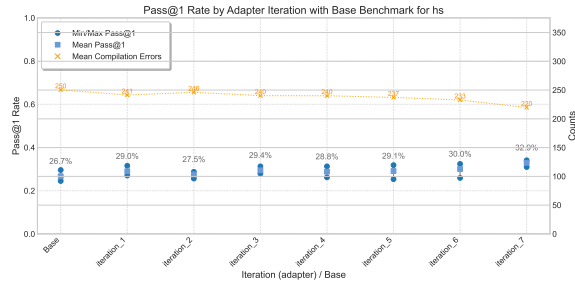


Figure 14: Alice's MBPP results. Averages over 16 trials.

H.2.3 Alice's SEQ vs SING Validated Generation Counts

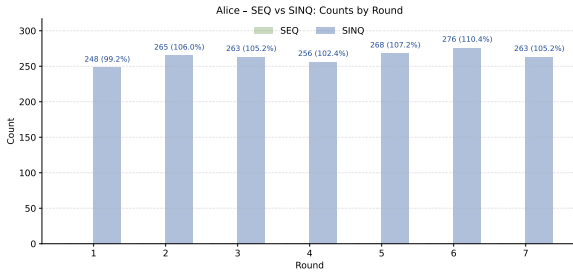


Figure 12: SEQ vs SING Validated Generation Counts.

H.3.3 Alice's SEQ vs SING Validated Generation Counts

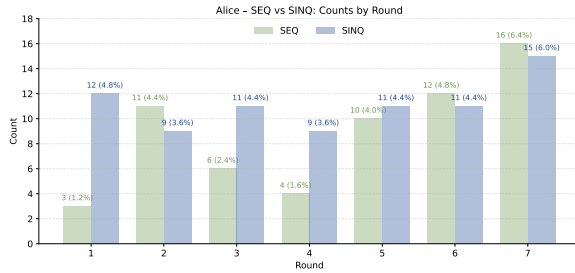


Figure 15: SEQ vs SING Validated Generation Counts.

H.3 Experiment E2

H.4 Experiment E3

H.3.1 Alice's HumanEval (Haskell) Performance

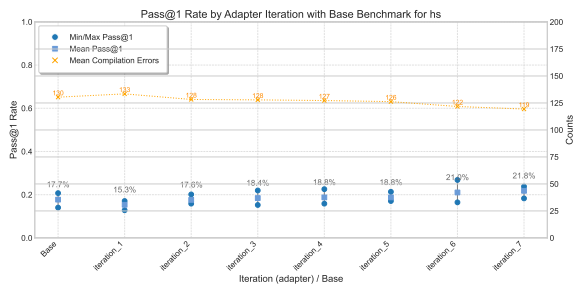


Figure 13: Alice's HumanEval results. Averages over 16 trials.

H.4.1 Alice's HumanEval (Haskell) Performance

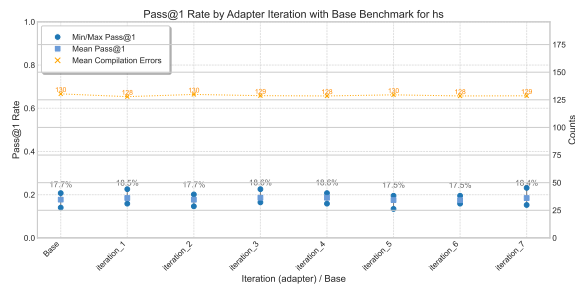


Figure 16: Alice's HumanEval results. Averages over 16 trials.

H.4.2 Alice's MBPP (Haskell) Performance

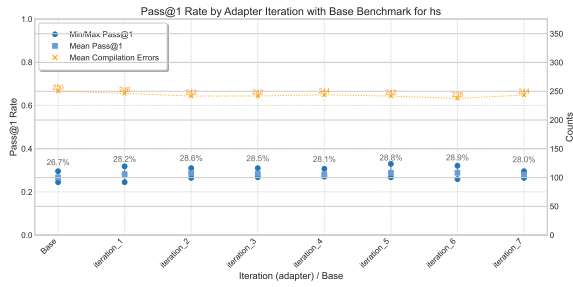


Figure 17: Alice's MBPP results. Averages over 16 trials.

H.4.3 Alice's SEQ vs SING Validated Generation Counts

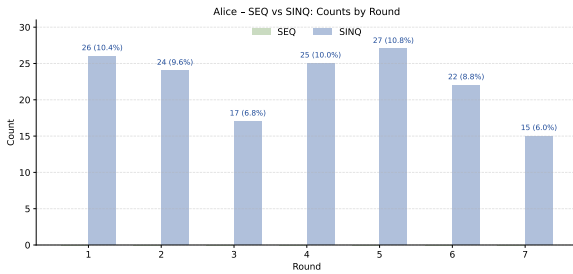


Figure 18: SEQ vs SING Validated Generation Counts.

H.4.4 E₂ vs E₃ Validated Generation Counts

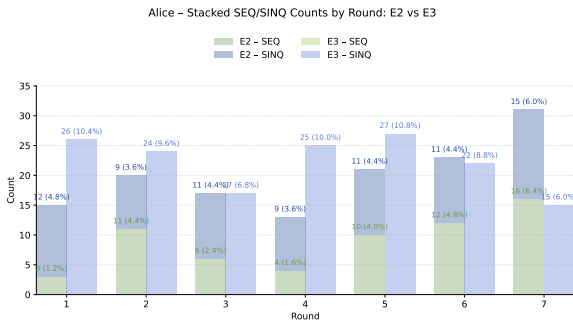


Figure 19: Validated Generation Counts of both experiment, indicating the effort of controlling P shown in Appendix F in order to achieve an equal number of verified generation.