

# The Verifier as Curriculum. VHG and the Third Role.

## The Verifier as Curriculum. VHG and the Third Role.

### A Daily Field Note on Three-Party Self-Play and Curriculum Construction

*Manu Bhardwaj. ifitsmanu.com. 10 May 2026. Last updated 10 May 2026. Version 1.0. Field Notes #5.*

Cite this article. Research index. Companion. The Asymmetry. Series origin. The Cost of Being Right.

**Daily field note.** First piece in the daily-review cadence. One fresh paper or post in the verification-economics or inference-economics wedge from the previous seven days, decomposed against what it shows and what it does not. Lighter scope than Field Notes #1 to #4. Same voice.

### What it claims

Lai, Feng, Teh, and Miao (2026), posted to arXiv on 8 May 2026, introduce VHG. A “verifier-enhanced hard problem generation framework built upon three-party self-play.” The setter generates problem-and-reference pairs. The solver attempts the problem. An independent verifier gates whether the setter’s problem counts as valid before the solver’s difficulty signal is applied. The setter’s reward becomes the indicator that the verifier accepts the problem times one minus the solver’s accuracy on it. Invalid problems get zero reward regardless of how hard they are. Reward hacking by emitting unanswerable or ill-defined problems no longer pays.

The headline numbers. On indefinite-integral generation evaluated against a Qwen3-4B-Base solver, VHG raises pass@1 from 52.5% to 69.4% on the Qualifier set, from 28.8% to 45.4% on the Competition set, and from 43.3% to 64.7% on the Stress Test set. On general mathematical reasoning, evaluated across five standard benchmarks (MATH, AMC, Minerva, Olympiad, AIME 2024 to 2026), pass@1 lifts from 56.8% to 69.0%. The lifts are larger for stronger solvers. Qwen3-14B reaches 49.23% pass@1 on integration and 41.50% on general math.

### How it argues

VHG instantiates two verifier variants. A *Hard* verifier built on SymPy that mechanically checks the validity of indefinite integrals. A *Soft* verifier instantiated as an LLM-as-judge for general mathematics. The Hard variant is the cleaner experimental probe. SymPy can decide “is this integral well-posed and does the setter’s reference solution differentiate back to the integrand” without ambiguity, so reward hacking through invalid integrals is foreclosed by construction.

The structural argument is that the conventional setter-and-solver duality is unstable. Without a

third party, the setter learns to maximize solver difficulty by generating problems that are hard because they are broken, not hard because they are deep. This is the standard reward-hacking failure mode in self-play with an LLM-graded reward. The verifier is the load-bearing constraint that re-aligns difficulty with validity. The reward function is multiplicative rather than additive, so a zero from the verifier kills the reward regardless of how much difficulty signal the solver supplies.

## What is interesting

The interesting structural property is that VHG extends the verifier’s economic role from a two-place job to a three-place job.

In Field Notes #2 and #3, the verifier had two production roles. Inference-time gating, where it acts on cost-correct’s denominator at decode time. And training-time reward function, where it acts as the RLVR signal that aligns the policy under verifiable rewards. VHG names a third role. Training-data curator. The verifier decides which problems enter the training distribution at all.

This third role is structurally distinct because it operates one level upstream of the training-time reward. RLVR optimizes the policy against a fixed verifier on a fixed problem distribution. VHG uses the verifier to construct the problem distribution itself. The same verifier artifact now governs three points in the production lifecycle: data, reward, decode.

A second observation. The Hard verifier outperforms the Soft verifier across the metrics where both are measurable, in the same direction predicted by the -asymmetry analysis in Field Notes #3. Verifier quality dominates other levers, and verifier quality is bounded by what the verifier can mechanically decide. The cleanest verifiers are domain-narrow, including SymPy on integrals, type-checkers on code, and unit tests on functions. The breadth-or-depth tradeoff shows up here as a quality-or-coverage tradeoff. Soft verifiers cover more ground with weaker guarantees. Hard verifiers cover narrow ground with stronger ones. The term in *Cost-correct* picks up the variance in either direction.

## What is missing

The paper reports lifts but several quantities that would make the framework directly comparable to the *Cost-correct* decomposition are not surfaced.

First. No explicit token-cost or compute-cost accounting at the data-generation step. VHG runs three models at training-data-generation time. The setter, the verifier, and the solver. The pass@1 lifts are reported against a fixed Qwen3 family and a fixed problem budget, but not against a fixed compute budget. Whether the lift survives at iso-FLOP across the full pipeline is the binding question for production adoption, and the paper does not answer it.

Second. The verifier-quality tradeoff is acknowledged but not quantified. The authors note that the Soft verifier “can still accept subtle errors, underspecified problems, or reward-hacking artifacts.” How often it does so, and how that propagates into solver quality drift over multiple self-play rounds, is left for future work.

Third. The empirical comparisons use one model family. The Qwen3 line is a strong open base, but not the only one. Whether the lifts transfer to other open bases or to proprietary frontier models is an open question, and the cost-economics framing would want it answered before generalizing the result.

## Why it matters now

Two reasons. The first is a curriculum-cost reason. Synthetic problem generation is becoming a binding input to frontier post-training as the natural-data math corpus gets exhausted. RL with verifiable rewards at scale needs more verifiable problems than humans produce, and self-play has been the obvious lever, with reward hacking the obvious failure mode. VHG offers a structural fix that does not require either manual curation or a stronger reward model. It requires only a tighter verifier.

The second is a regime-fit reason. The *Cost-correct* equation from Field Notes #2 puts  $\gamma$  in the denominator. Field Notes #3 shows that  $\gamma$ -engineering is the highest-leverage place to spend a marginal engineering dollar in the typical operating regime. VHG identifies a fourth way the same engineering dollar can land. Not as an inference-time gate. Not as an RL reward function. Not as a process reward model. As a curriculum constructor whose output is the training distribution itself. The three-party game is the production architecture that closes the loop. Generator, verifier, curriculum, all sharing one verifier artifact.

The cleaner the verifier gets, the more of the lifecycle it touches. That is the empirical pattern this paper extends and the analytical lever it adds to the framework.

---

## Source

- Lai, Y., Feng, J., Teh, Y. W., and Miao, N. *Verifier-Backed Hard Problem Generation for Mathematical Reasoning*. arXiv:2605.06660, 8 May 2026.

## Related field notes

- Field Notes #2. The Cost of Being Right. Verification Economics in 2026.
- Field Notes #3. The  $\gamma$  Asymmetry. Why Verifiers Can Be Smaller Than Generators.

---

## Cite this article

---

Research index. Home.