

Verifiable by Design: Synthetic Data, Reinforcement Learning, and Where AI Improves

Luís Seabra

Protegrity · 2026/04/14

Abstract

In 2026, AI is transforming software engineering substantially, and in ways few predicted. Understanding why tells us more about the current state of field than any benchmark score or lab announcement. The answer turns out to be less about what the models can do and more about the kind of data that can train them, and specifically whether that data can be verified.

This paper argues that the dramatic progress visible in AI coding agents is best understood through a data problem: not just generating synthetic training examples at scale, which any capable model can do cheaply, but knowing which of those examples are actually correct. When you can grade outputs as correct or incorrect without human judgment, synthetic data becomes trustworthy and reinforcement learning becomes productive. This ability is what we call the verification function, and its availability is what determines where synthetic data pipelines can be built and iterated, and where they cannot. Software engineering uniquely satisfies this condition through executable test suites that return a clean binary signal. Most other domains do not, and the more productive question for 2026 is not whether AGI is approaching, but where else a reliable verification function might be engineered, and what synthetic data pipeline could follow.

I. Every Past Has Its Future — And Its Data

There is a pattern to how we imagine artificial intelligence, where we always project what we expect it to be in a different direction than what eventually happens. Every generation of AI research has produced a vision of the future, and, as so often happens with science fiction, it does not age well in its specifics. The influence tends not to flow from the prediction but from wherever the previous wave happened to find solid ground — and solid data.

In the 1980s, the dominant paradigm was symbolic AI and expert systems, where hand-crafted rules would encode the knowledge of specialists, promising machines that could reason from knowledge, making doctors or lawyers obsolete. The data of that era was rules: these were curated, brittle, and expensive to maintain. Soon these systems hit walls that no amount of rule-writing could scale. In the early 2010s, the explosion of deep convolutional neural networks, training on ImageNet, running on GPUs, redirected the field toward perception layers with an initial Computer Vision focus. The data of that era was labelled images, assembled at industrial scale through crowdsourcing and annotation pipelines. AI was now mostly imagined on how machines would see and recognise the world, the first step of understanding it. Radiologists would be the first profession to go. Self-driving vehicles were the flagship example with Level 5 autonomy, repeatedly, a few years away.

Between 2017 and 2020, the transformer architecture with its attention mechanism reshuffled projections again. BERT, GPT and their successors demonstrated that language could be modelled at scale allowing useful, general-purpose next-token predictors to build natural language processing pipelines. The data of that era was the Internet itself: billions of documents, scraped, filtered, and fed into training runs of unprecedented scale. Simply by being prompted in the right way, models were able to do zero-shot and few-shot adaptation to tasks they had not been explicitly trained on. The predictions adjusted accordingly: universal knowledge assistants, AI tutors, agents managing your calendar and finances. Then in late 2022 ChatGPT brought these capabilities to a mass audience in a single product, and the excitement ramped off the charts. It was genuinely useful for drafting, researching, explaining. This cross-domain usefulness and its fluency signaled to many that something qualitatively new had arrived. The excitement cycle took off: Artificial General Intelligence, or AGI, was weeks away, or months, or perhaps already here depending on whom you asked. The hype around AI agents as personal companions, systems that knew you and would act for you autonomously in the world, reached a rhetorical peak through 2024, when Gartner placed autonomous agents squarely at the top of its Hype Cycle for Artificial Intelligence [Gartner2024].

It is now 2026. While prominent voices in the field claim AGI is imminent or already functionally here, other equally prominent voices remain openly skeptical, arguing the benchmarks are saturating faster than the capabilities are generalizing. In practice, and despite the positive spin, personal AI companions have not yet arrived, save for a handful of DIY enthusiasts and entrepreneurs wiring something together themselves. AI doctors, including radiologists, and lawyers remain tools, not practitioners. Self-driving in a scalable way is proving harder and more expensive than the predictions allowed for. And yet the field AI is actually transforming is the one closest to home for its creators: software engineering.

This paper asks why, and the answer turns out to be less about what the models can do and more about the kind of data that can train them.

II. On Intelligence

Before pondering whether AI is generally intelligent, it is worth framing what intelligence is, a question that has occupied philosophers and cognitive scientists for decades without producing a consensus, and which I do not plan to address here.

Textbook definitions involve some combination of learning from experience, applying knowledge to novel situations, reasoning under uncertainty, and, crucially, the capacity to pursue goals in a dynamic, complex environment. A useful distinction, drawn by psychologists like Kahneman, is between fast, automatic cognitive processes, pattern recognition, usual responses, instinct, and slower, deliberate reasoning [Kahneman2011]. A large proportion of what we call thinking, in everyday human life, does not actually involve what we would want to call intelligence. Looking up a word, recalling a social convention, responding to a question with a rehearsed answer: these draw on prior learning and serve adaptive functions, but they are not what defines intelligence as a distinctively human capability. They are closer to sophisticated, learned instinct than to generative reasoning.

This distinction matters when evaluating current AI systems. Auto-regressive models learn a probability distribution over training text and produce the statistically likely continuation of whatever has been said so far. This description understates what they can visibly accomplish: extract concepts, follow novel instructions, and generalize across situations they were never explicitly trained on. The outputs are too coherent, too contextually sensitive, too useful for a simple dismissal. And yet they are, in essence, a very fast and very well-calibrated instinct engine running on the accumulated residue of human reasoning, compressed and made retrievable.

Impressive is not the same as intelligent, a distinction which I do not think is merely cosmetic. What LLM-powered AI has learned is better described as the shape of intelligent output, not the process that produces it. Their first instinct is to answer, comply, and execute. They do not question the premise, interrogate the context, or push back and try to change the world around them. They are machines that answer questions, not sentient beings that ask them.

III. The Data Problem and why Synthetic works for Software

The reason these engines are succeeding in software is not that they are intelligent. It is that software gives them something most domains cannot: a way to know, automatically and without human judgment, when they got it right. To answer why, start with data.

Something remarkable is happening in software engineering in 2026. In February 2025, Andrej Karpathy coined the term 'vibe coding' [Karpathy2025] to describe fully delegating code generation to an AI, completely sidestepping any human reviewing. It was originally conceived for throwaway weekend projects and widely dismissed as a toy, something fit for hobby apps and demos but with no place in production systems where real consequences follow from bad code. Then something unexpected happened: it worked. Not perfectly, not safely by traditional standards, but well enough that developers started using it for real work, then teams, then companies. Within a year the conversation had moved from whether vibe coding was serious to how to scale it. Coding agents are now decomposing requirements, making architectural decisions, writing tests, catching their own errors, and shipping working software, all with minimal or no human supervision. The term 'agentic engineering' is on the rise [Karpathy2026] to support this movement.

What emerges looks like competence, something which may be uncomfortable for those who have always thought of software engineering as a craft. The question I find revealing is not whether this is a signal of approaching general intelligence, but what it is about software engineering that made it the first domain to be transformed.

The Data Wall

The scaling era of internet-scraped text is approaching its limits. The stock of high-quality human-written text on the internet is finite, and frontier models have consumed most of what is worth consuming. Projections suggest that at current training scales, models could exhaust the available supply of unique, useful web text within this decade [VillalobosEtAl2024]. The returns from scaling pre-training data are yielding diminishing returns, and the field has been looking for alternatives. Synthetic data, text

generated by models or automated pipelines rather than by humans, is the most promising alternative, but it comes with a catch.

Training models on their own outputs, if done naively, leads to what is called model collapse: a gradual narrowing of the output distribution, where the model becomes increasingly confident about an increasingly restricted range of responses [ShumailovEtAl2024]. The outputs grow stale, rare knowledge fades, and the model converges on the average. You end up with more data and a worse model.

Generating useful synthetic data requires more than prompting a model and collecting outputs. In practice, effective pipelines seed generation with real examples to anchor diversity, use structured-template based prompting strategies or dedicated fine-tuned models to reach the edges of the distribution, and apply filters to ensure rare but important cases are not crowded out by the modal majority. The collapse risk is highest precisely here: a model asked to generate its own training data will tend toward confident, common outputs and systematically underrepresent the long tail. Verification addresses this not by constraining what gets generated, but by selecting or labeling what gets used for training. Rare correct examples survive the verification filter just as common ones do, which preserves distributional coverage in a way that generation alone cannot guarantee.

The response from the research community has been telling. HuggingFace's FineWeb-Edu filters a 15 trillion token web corpus down to 1.3 trillion tokens of high-quality content, using classifiers trained on synthetic annotations, following the same approach used by Llama 3 and Phi-3 in their non-public pipelines [LozhkovEtAl2024]. Even the curation of raw pre-training data now depends on synthetic data pipelines to identify what is worth keeping.

Synthetic data generation is cheap but verification is hard. Any capable LLM can produce training examples at scale, for almost any domain, at negligible cost. Knowing which of those examples are actually correct, and worth learning from, without a human in the loop, is where the difficulty lies. And knowing which domains even allow for this kind of automatic grading is, arguably, the more important question.

The Verification Function

If you have a function that can grade the output of your generative process as correct or incorrect without human judgment, you have an infinite supply of labelled training data. You generate a candidate via some prompting or a crafted model, you run it through the verification function, you use the result as a training signal. It is the verification function that imparts knowledge to the pipeline, not the generator, which is blind, but the grader, which knows what good looks like.

This is precisely what modern reinforcement learning training pipelines are designed to exploit. Techniques like Proximal Policy Optimization (PPO) and Group Relative Policy Optimization (GRPO) allow models to improve not by memorizing more examples but by being rewarded for solutions that pass and penalized for solutions that fail [Schulman2017; ShaoEtAl2024]. Over many iterations, the model generalizes across the whole class of problems, not because it memorized solutions, but because it was trained on an effectively unlimited supply of freshly generated problems, each verified automatically. The verification function is what makes this possible: it grades each candidate solution, turning the

output of a code generator into a labelled training example at zero marginal cost. Without it, there is no RL loop. Without an RL loop, there is no systematic improvement beyond what the pre-training data already contains.

This also means verification becomes the pacing constraint for the overall pipeline, not generation. A recent 14-billion-parameter coding model trained on 32 H100 GPUs required over 41,000 sandboxed code executions per training step, demanding a dedicated CPU verification service capable of running more than 100 concurrent environments simultaneously [LuoEtAl2025]. The model learned faster when the verifier could keep up, not when the generator ran faster.

Why Software is Special

Software development uniquely provides a verification function of extraordinary quality. A software task has a specification — requirements, a user story, a bug report. It has an implementation — the code. And it has a test suite: a set of executable checks that determine, without ambiguity, whether the specification is implemented. Thousands of different problems can be constructed programmatically, run the tests, and get a clean binary signal back [Jimenez2024]. No taste, no interpretation, no cultural context. Just pass or fail.

This is not just a convenient property of coding tasks. It is, in a technical sense, a verification function, enabling exactly the kind of synthetic data generation pipeline that the field was searching for. You generate a problem statement programmatically, produce thousands of candidate solutions, and run the test suite on each. The passing ones become positive training examples; the failing ones become negative ones. The synthetic dataset is unlimited because the problem generator is unlimited. The labels are free because the test runner is free. The training signal is clean because the tests are deterministic.

Certainly, there are other factors at play with software. The sheer volume of publicly available code provides pre-training data at a scale few other domains can match. Programming languages have formal grammars, unambiguous and machine-parseable in ways that natural language is not. Code is compositional, built from discrete reusable blocks layered across levels of abstraction. Tasks decompose naturally into subtasks, each independently verifiable. These properties do not replace the verification function. They are why, in software, it works so well.

The results are visible in benchmarks. SWE-bench, the standard evaluation for coding agents on real-world GitHub issues, has had to be repeatedly revised and hardened as agents keep saturating it [Jimenez2024]. That is not a problem with the benchmark but rather evidence that the pipeline is working. What is happening in software engineering in 2026 is not evidence of the AI discipline approaching general intelligence methods. It is, I argue, that software engineering is the domain where synthetic data and reinforcement learning are working in a tight loop to generate, verify and train.

From Model to System

A base model alone, no matter how powerful, is not sufficient to achieve this level of proficiency in software engineering tasks. The coding agent harness built around it is a runtime, inference orchestration loop that equips the model with tools, tracks state and feeds the output of each step back as the input of the next. We observe echoes

of the synthetic data pipeline structure: constrain the process, verify the output, retry on failure. The harness works for exactly the same reason the training loop does. Every step produces a signal that can be checked. The current generation of tools, Claude Code and its competitors, shed the IDE scaffolding entirely. The user expresses intent in plain text; the agent researches, plans, and executes from there. The system as a whole produces something that looks like engineering. But it is, at bottom, a very sophisticated instinct machine held accountable by automated tests.

The frameworks being developed on top of coding agents lean on this insight. The Ralph pattern [Huntley2025] runs a coding agent autonomously in a loop against a product requirements document, iterating until all tasks pass their tests. Ralph does not learn between iterations. He is not intelligent, but rather persistent, and persistence inside a verifiable loop turns out to be valuable. Gastown [Yegge2026] scales this further, coordinating twenty to thirty parallel agents on the same codebase. Serial or parallel, simple or orchestrated, the architecture is the same. What the harness cannot do is care about the difference between code that works and code that is secure, efficient, maintainable, or designed for a future nobody has yet specified. While verification confirms function, this judgment remains human.

IV. The Valley of the Unverifiable

AI has not had the same impact on music, film, or literature that it has had on software engineering, and the gap is not for lack of trying. The financial incentives are enormous. The models have properties that superficially mimic creativity: temperature and sampling introduce controlled randomness, a mechanical proxy for artistic spontaneity. It is tempting to call this a matter of timing. I think it runs deeper.

Systems like Suno and Udio can generate music in any style you care to name. Midjourney and Stable Diffusion produce striking images. Sora generates video that, in short clips, can be genuinely impressive. Since the initial wave of excitement, these tools have settled into a more modest place in the cultural landscape, useful for mockups and prototyping, but not much else. They are not producing work that people form lasting relationships with. The outputs are technically accomplished but aesthetically inert. These systems differ substantially under the hood — diffusion models, transformer hybrids, autoregressive generators — but that is not where the explanation lies. Generating synthetic training data for music or images is not the bottleneck; that too is trivially easy. What does not exist is a way to label it. What they all share is the absence of a verification function capable of guiding improvement toward the thing that actually matters.

There is something about this aesthetic inertness that recalls what might be called an uncanny valley effect, though the discomfort operates differently here than in the original formulation. It is not that the outputs look almost-but-not-quite human. They feel like the average of everything rather than the specificity of anything. Great art, even great popular art, is recognizable as coming from someone, as the product of a particular sensibility encountering the world in a particular way. AI-generated creative content is, by construction, the synthesis of a distribution. Humans are quite sensitive to this, and the absence of a voice is striking once the initial novelty wears off.

Through the lens of synthetic data, this failure is structural, not contingent. Generation is not the hard part. We can produce synthetic music, images, and video at scale. The problem is that we cannot tell which of it is good, and once again, the bottleneck is verification. There is no function that grades a piece of music as good or bad in the way a test suite grades code. You can generate synthetic training examples but labels require taste, and taste is precisely what cannot be formalized. Even human labels disagree in ways that make them a weak training signal. The RL loop has nothing to grip.

Consider a thought experiment: what if we tried to build a verification function for artistic quality? Google's MusicRL attempted exactly this, fine-tuning a music generation model using reinforcement learning on human preference signals [Cideron2024]. What it produced was better alignment with average preferences, not an artistic voice. Mimicry, not music. You cannot write a test suite for the property that actually matters, the property of producing work that feels like it came from someone rather than something. The gradient does not exist as a computable function.

This feels like a genuine ceiling of the current paradigm, and not one that more compute will move. The bottleneck is not in capability, but rather the inability to formalize what good actually means in these domains. For domains where human judgment about quality is subjective, diffuse, or culturally constructed, the synthetic data pipeline that has been so productive in software engineering simply has nothing to anchor to. The same pattern is currently observed in science: AlphaFold succeeded at protein structure prediction because experimental ground truth provides a verification function. Drug efficacy does not have one. The feedback signal is a clinical trial, which takes years and cannot be generated synthetically.

Within software engineering, the verification function weakens as you move up the abstraction hierarchy, and with it so does the usefulness of synthetic data. From passing tests to good code, from good code to right architecture, from right architecture to correct product, each step is harder to verify and therefore harder to generate useful training signal from. Deciding what to build, anticipating what users will need before they know it themselves, knowing when to push back on the requirements entirely: these are the things that distinguish an experienced engineer from a code generator, and no synthetic pipeline knows how to label them yet. And the harder the judgment call, the more costly verification becomes to run at scale, until the compute required exceeds what any training pipeline can practically sustain.

V. What the Future Tells Us

AI is advancing fast, and what seemed implausible a year ago is routine today. But capability and speed are not the same as intelligence, and that distinction is worth holding onto. The word intelligence is thrown around too freely. What we have right now are systems extraordinarily good at generalizing from training samples within the category of problems where the right answer can be computed, checked, and

rewarded. This is a significant and economically consequential subset of cognitive labor, but it is not an intelligence signal.

The discourse tends toward an AGI-or-nothing binary that obscures more than it reveals. A more productive question, and the one this paper has been building toward, is: where else can we engineer a synthetic data pipeline grounded in a reliable verification function? Where we can construct a reliable verification function, synthetic data pipelines can be built, iterated, and improved. Where we cannot, we are either looking at domains that require human judgment, or at genuinely open scientific territory. That boundary is not fixed. It can be pushed by better tooling, better problem decomposition, and occasionally by finding that a domain has more verifiable structure than it first appeared. But it cannot be wished away. This maps loosely onto a familiar distinction in computer science: problems where solutions are hard to find but easy to check are precisely the ones where automated search can make progress. Where verification itself is computationally hard or requires human judgment, the search has nothing to guide it.

Some domains look more promising than others for applying verifiable synthetic data pipelines. Mathematics is an early indicator: proof assistants like Lean already provide formal verification functions and benchmark results are following. Formal logic and structured puzzles share the same properties [Chen2025]. Mathematics may also have a hidden advantage in that verification is fast. Checking a numerical answer costs under a millisecond; code execution in a sandbox costs one to ten seconds. This would interestingly imply the tractability of a domain would depend not just on whether a verification function exists, but on how cheap it is to run at scale. Structured information retrieval is another case, where whether a query returns the correct result is as verifiable as whether code passes a test. Compliance and regulatory checking may matter most economically. Wherever rules can be formalized precisely enough to encode, the pipeline becomes tractable: synthetic transactions, documents, or agent actions can be generated and checked against the rule automatically, with no human in the loop. Legal reasoning and protection policy enforcement share this structure wherever statutes and precedents are precise enough to function as ground truth.

If this approach has any predictive value, the domains to watch are those where verifiability and economic pressure coincide. Legal reasoning, regulatory compliance, and software access governance all have both, and the incentives to build the pipeline are already there. Where a reliable verification function can be engineered, the benchmark saturation pattern we saw in software will repeat.

Whatever lies beyond the current architectural paradigm on the road to something akin to general intelligence, my instinct is that the path runs beyond the transformer, perhaps toward something more like world models [LeCun2022]. What is certain is that whatever architectures emerge, they will need data, and the question of how to generate and verify it at scale will not go away. If anything, it becomes harder as the problems grow more ambitious and the domains less structured. Synthetic data is not a transitional workaround, it is the core challenge.

We may learn more about intelligence by asking empirically, domain by domain, where the verification function can be built and where it cannot. Any progress toward something like general intelligence will have to be grounded in good, well-verified data.

References

- LeCun, Yann. "A Path Towards Autonomous Machine Intelligence." OpenReview, June 2022. <https://openreview.net/forum?id=BZ5a1r-kVsf>
- Karpathy, Andrej. Post on X, February 2, 2025. <https://x.com/karpathy/status/1886192184808149383>
- Karpathy, Andrej. Post on X, February 2, 2026. <https://x.com/karpathy/status/2019137879310836075>
- Yegge, Steve. "Welcome to Gas Town." Medium, January 2026.
- Huntley, Geoffrey. "Ralph is a Bash Loop." Claude Code community methodology, 2025.
- Gartner. "Hype Cycle for Artificial Intelligence, 2024." Gartner Research, 2024.
- Jimenez, Carlos E., et al. "SWE-bench: Can Language Models Resolve Real-World GitHub Issues?" ICLR, 2024. arXiv:2310.06770.
- Cideron, Geoffrey, et al. "MusicRL: Aligning Music Generation to Human Preferences." ICML, 2024. arXiv:2402.04229.
- Shumailov, Ilya, et al. "AI Models Collapse When Trained on Recursively Generated Data." Nature, 2024.
- Villalobos, Pablo, et al. "Will We Run Out of Data? Limits of LLM Scaling Based on Human-Generated Data." arXiv:2211.04325, 2024 update.
- Lozhkov, Anton, et al. "FineWeb-Edu: the Finest Collection of Educational Content." HuggingFace, 2024. <https://huggingface.co/datasets/HuggingFaceFW/fineweb-edu>
- Kahneman, Daniel. Thinking, Fast and Slow. Farrar, Straus and Giroux, 2011.
- Schulman, John, et al. "Proximal Policy Optimization Algorithms." arXiv:1707.06347, 2017.
- Shao, Zhihong, et al. "DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models." arXiv:2402.03300, 2024.
- Chen, et al. "Enigmata: Scaling Logical Reasoning in Large Language Models with Synthetic Verifiable Puzzles." arXiv, 2025.
- Luo, Michael, et al. "DeepCoder: A Fully Open-Source 14B Coder at O3-mini Level." Together AI Blog, April 2025. <https://www.together.ai/blog/deepcoder> [LuoEtAI2025]