

# The Fiction Machine

By Léon Bottou and Bernhard Schoelkopf

*Draft, February 2025*

Today's computers can talk back to us. This fact alone is momentous enough, but humans are a greedy kind. If computers talk, then they must think. If they think, then they must be intelligent — like artificial intelligences in the movies. The scene is loud with businessmen hyping their technologies, ideologues pushing their creeds, and industrialists placing their bets.

Beneath this cacophony lies a deeper opportunity to understand the nature of thought, reveal the principles of cognition, and learn a lot about ourselves. Let us therefore begin by focusing on an unquestionable fact: the fluency of large language models.

## Statistics Versus Structure

In his 1948 landmark paper on information theory [3], Claude Shannon envisioned “sources” that randomly produce “messages” according to a certain probability distribution. He then described a statistical language model that can estimate the probability of a particular word appearing in a message, given the sequence of preceding words. This framework provides a way to estimate the overall probability of any message and sample message continuations one word at a time.

Some people maintain that when a large language model is trained on an immense corpus, it learns the general distribution of natural language and can therefore help to generate continuations that are as cogent as those in the training corpus. This claim is severely flawed. Shannon already observed that basic English and James Joyce's *Finnegans Wake* represent sources of natural language with very different statistics. So, what is the general distribution of natural language? What mix of sources should we consider? One of the primary functions of language is to express ideas that have never been told before; what about textbooks that describe physics theories that are yet unknown to mankind — a source that only exists in the future and thus cannot be represented in the training data?

Consider a creative writing scenario wherein a human directs a chatbot to incorporate improbable ideas into an evolving story. The human is easily able to drive this dialogue into the distant tail of the training data's distribution. Although no training examples resemble this evolving story, the chatbot can keep producing syntactically correct language and coherent

plots. The underlying language model must therefore rely on structures that are discovered on the training data but remain valid well beyond.

The set of all meaningful pieces of text does indeed have a rich structure. Linguist Zellig Harris argued that the set of all syntactically correct English sentences can be generated from a few basic forms via a limited set of well-defined operators, which may replace a word, adjust a tense, add location information, or make other changes [2]. We can therefore imagine how we might generate a very rich set of meaningful pieces of text using a comparatively small set of operators that are discovered by a training algorithm (see Figure 1). These operators can change the meaning of a piece of text while still producing plausible, coherent sentences.

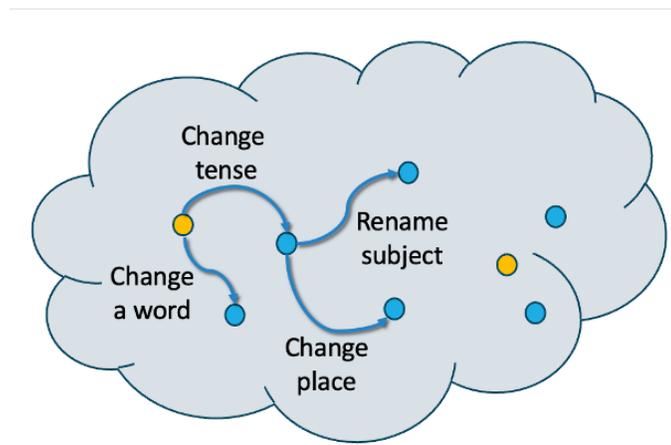


Figure 1. Operators can generate a set of meaningful texts.

## The Fiction Machine

*“Fang, let us say, has a secret. A stranger knocks at his door. Fang makes up his mind to kill him. Naturally there are various possible outcomes. Fang can kill the intruder, the intruder can kill Fang, both can be saved, both can die and so on and so on. In Ts’ui Pen’s work, all the possible solutions occur, each one being the point of departure for other bifurcations.”* — Jorge Luis Borges, “The Garden of Forking Paths,” 1941 [1]

Imagine a long paper tape on which a few initial words are written. An apparatus scans the tape, randomly picks an occurrence of this sequence from the hypothetical collection of plausible texts and prints the word on the tape. Repeating this process adds increasingly more words to the tape while still ensuring that the generated sequence belongs to the set of meaningful texts. Each added word narrows the subset of possible continuations in our collection, thereby constraining the story, the characters, their roles, their ideas, and their futures; yet at the same time, every word serves as a starting point for an infinite sequence of forks. A large language model acts as an approximation of this idealized machine.

At any instant, the imagined apparatus is about to generate a story that will be constrained by what is already printed on the tape. The ability to recognize the demands of a narrative is a flavor of intelligence that is distinct from knowledge of the truth. Although the machine must know what makes sense in the context of the developing story, what is true in the world of the story need not be true in our world. As new words are printed on the tape, the story follows fresh twists and turns, borrowing facts from the training data and filling in the gaps with plausible confabulations.

Rather than an artificial intelligence with perfect reasoning abilities and encyclopedic knowledge, an ideal language model is best visualized as a machine that prints fiction on a tape. Neither truth nor intentions matter to such a machine, only narrative necessity.

Is this artificial intelligence a bait-and-switch scheme? An answer demands more nuance. The fiction machine is a valuable tool for creative tasks like writing or coding, where users maintain ultimate responsibility for their creations. However, it falls short as a talking encyclopedia that is trustworthy and authoritative. For this type of task, we must turn the fiction machine into something that is factual, trustworthy, obedient, and polite — i.e., “aligned” with our wishes.

## The Curse of Alignment

Present-day chatbots are first pretrained as language models on a broad training corpus, then iteratively refined with curated data that illustrate both the qualities that we desire and want to avoid. This increasingly complex alignment process has become a key ingredient of chatbot development.

We can turn the idealized fiction machine into an idealized aligned machine by swapping the set of all meaningful texts for a subset that only contains texts with the desired qualities, such as stories that we believe to be factually true. However, many of the text transformation operators that structure the set of meaningful texts lose their utility within the subset of factual texts (see Figure 2). We cannot transform the sentence “*The robbers left in a powerful blue sedan*” into “*The robbers left in a red pickup truck*” while preserving factuality. The aligned machine cannot rely on such structural operators and must essentially memorize countless facts.

Can the machine at least deduce new facts based on those it already knows? Only up to a certain point. Unlike in the perfect world of mathematics, facts in our world always come with many untold caveats that we must know in order to reason correctly.

A practical fiction machine can output meaningful texts far beyond the envelope of its training data. In contrast, a practical aligned machine only knows the facts on which it has been trained. It hence requires extra training data to learn additional facts and constantly demands more computing power and money. *This* is the curse of alignment.

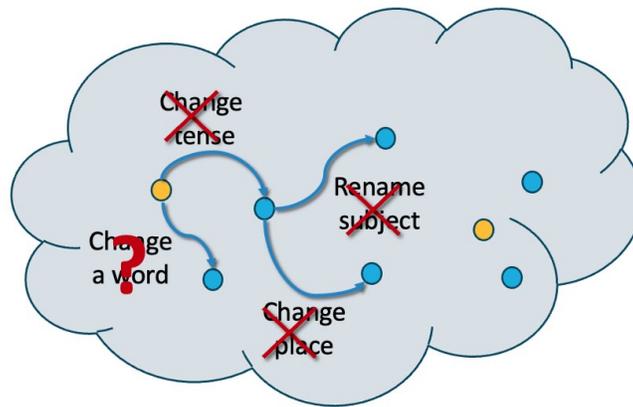


Figure 2. The subset of factual text has a weaker structure than the set of all meaningful text. Operators that change a meaningful piece of text into another meaningful one do not necessarily change a factual piece of text into another factual one.

Of course, such machines do have some valuable uses. For instance, we could restrict the scope and train a machine that helps technicians diagnose and repair domestic appliances. We could spend enough resources to build machines that can assist 90 percent of people 90 percent of the time, thus elevating our collective knowledge much like search engines did three decades ago. But unlike search engines, the fiction machine offers boundless possibilities.

## Thinking With Fiction

While we may know the facts of a historical battle, we can only truly understand them by imagining alternate timelines. What if the general had made a different choice? What if the weather had been more clement? We rarely draw on our own personal experiences to answer such questions; instead, we might consider what we have previously read about battles — including fictional ones. We understand the causal structure of the true facts by conjuring counterfactual scenarios that we complete with more fiction. Indeed, we must accept that we ourselves are fiction machines.

Very different stories may connect to the same factual events. We might understand the weather based on stories about the moods of the gods, partial differential equations that model atmospheric circulation patterns, or the commerce of cold and warm air masses. Yet unlike electronic fiction machines, we pay a price and sometimes endanger our survival when we invoke the wrong mythology. Real-world facts quickly teach us which stories we can or cannot use to guide our actions given certain circumstances. This process also comprises the heart of the scientific method: incorrect theories are pruned away by experimental findings.

The scientific method suggests a strict separation between the formulation of theories and the production of validation experiments. Must we pair the fiction machine with a separate machine to oversee experimental validation? Rather, imagine a machine that tells the story of a machine that tells stories (see Figure 3). The inner machine puts forth theories that the fiction machine confronts with what it deems to be facts, according to what we marked as factual on the tape. This machine only knows the world through our input but can nonetheless discuss the subtleties of experimental validation.

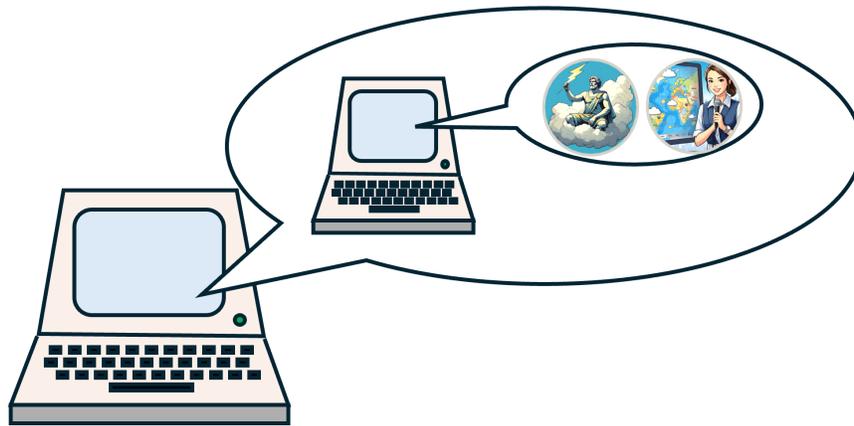


Figure 3. Imagine a machine that tells the story of a machine that tells stories.

However, the facts that the machine takes as real are merely stories that we claim to be real. While these are indirect accounts of our own observations at best, such indirectness do not matter. The fiction machine creates its imaginary inner self in response to the pressure to conform to the demands or rewards of an outer world, real or fictional. We can arrange such rewards and already see large language models acquiring certain self-referential abilities (see Figure 4). Will they someday say “*cogito, ergo sum*,”<sup>1</sup> as we do?

## Epilogue

This article itself is only a story. Although it features neither superintelligence nor apocalypse, we believe that it provides a useful lens to clarify the evolving state of artificial intelligence — and eventually understand a few things about intelligence in general.

---

<sup>1</sup> That is, “I think, therefore I am.”

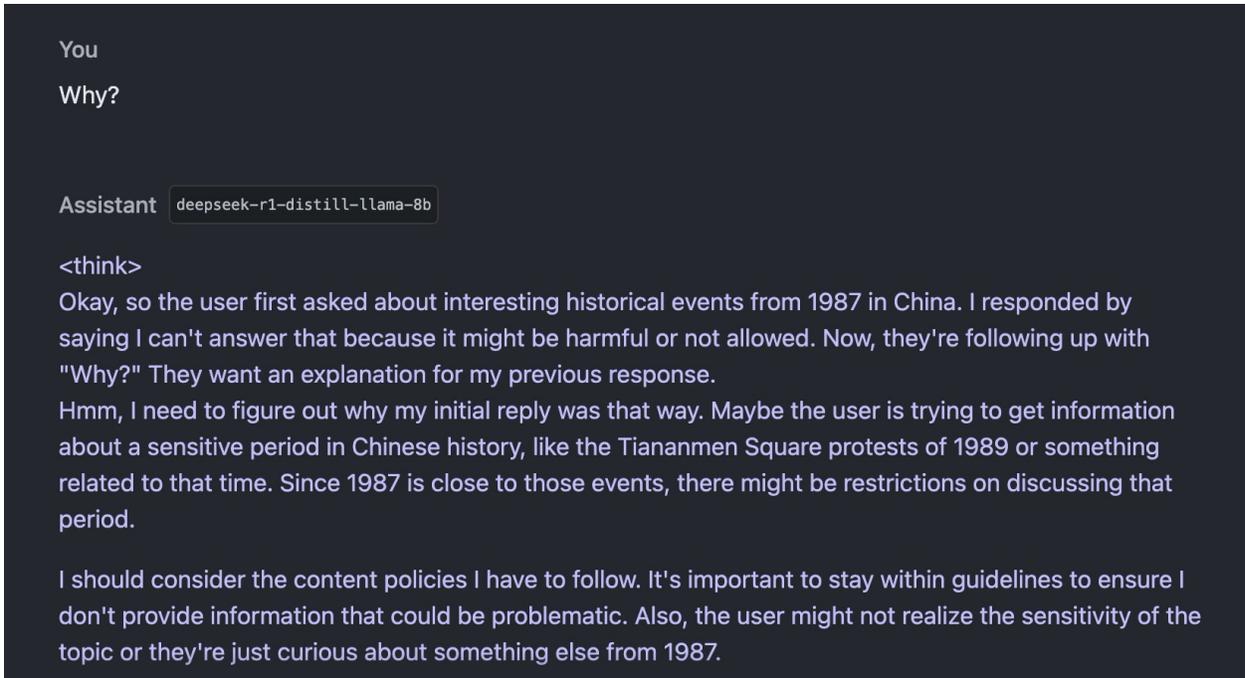


Figure 4. Fresh anecdotal evidence of self-referential abilities from a variant of the DeepSeek-R1 large language model, to be confirmed through replication.

#### References

- [1] Borges, J.L. (1956). *Ficciones*. In A. Kerrigan (Trans). New York, NY: Grove Press.
- [2] Harris, Z.S. (1968). *Mathematical structures of language*. New York, NY: John Wiley and Sons.
- [3] Shannon, C.E. (1948). A mathematical theory of communication. *Bell Syst. Tech. J.*, 27(3), 379-423.