

BEYOND POPULAR SCIENCE



DAVID H. SILVER



BEYOND POPULAR SCIENCE

David H. Silver

<https://www.openbookpublishers.com>

© 2026 David H. Silver



This work is licensed under the Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC 4.0). This license allows you to share, copy, distribute and transmit the text; to adapt the text for non-commercial purposes of the text providing attribution is made to the authors (but not in any way that suggests that they endorse you or your use of the work). Attribution should include the following information:

David H. Silver, *Beyond Popular Science*. Cambridge, UK: Open Book Publishers, 2026,
<https://doi.org/10.11647/OBP.0526>

Further details about CC BY-NC licenses are available at
<https://creativecommons.org/licenses/by-nc/4.0/>

Copyright and permissions for the reuse of many of the images included in this publication differ from the above. This information is provided in the captions and in the list of illustrations. Unless otherwise stated, figures are reproduced under the fair dealing principle. Every effort has been made to identify and contact copyright holders and any omission or error will be corrected if notification is made to the publisher.

All external links were active at the time of publication unless otherwise stated and have been archived via the Internet Archive Wayback Machine at
<https://archive.org/web>

Digital material and resources associated with this volume are available at
<https://doi.org/10.11647/OBP.0526#resources>

ISBN Paperback:	978-1-80511-877-0
ISBN Hardback:	978-1-80511-878-7
ISBN Digital (PDF):	978-1-80511-879-4
ISBN HTML:	978-1-80511-881-7
ISBN Digital ebook (epub):	978-1-80511-880-0
DOI:	10.11647/OBP.0526

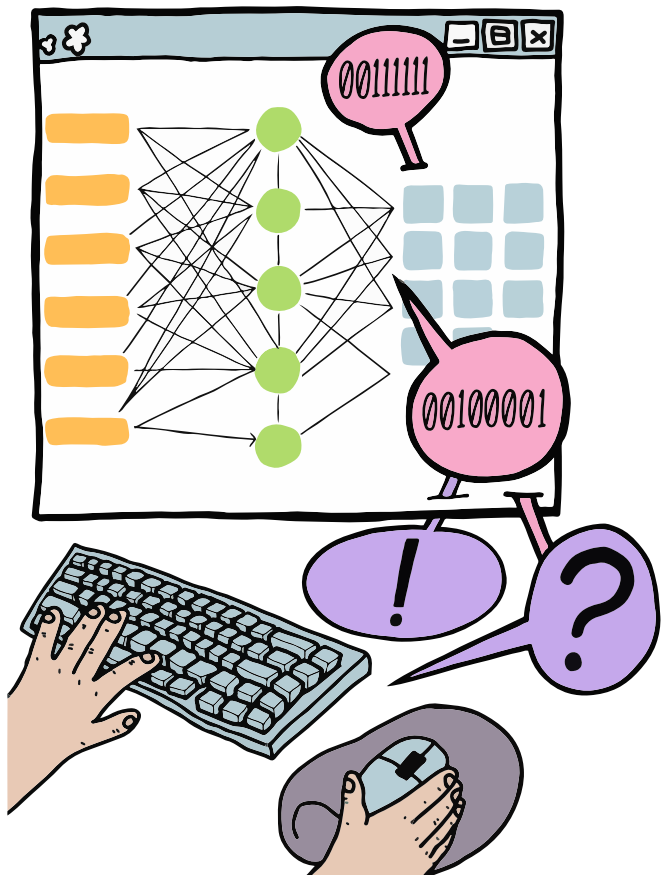
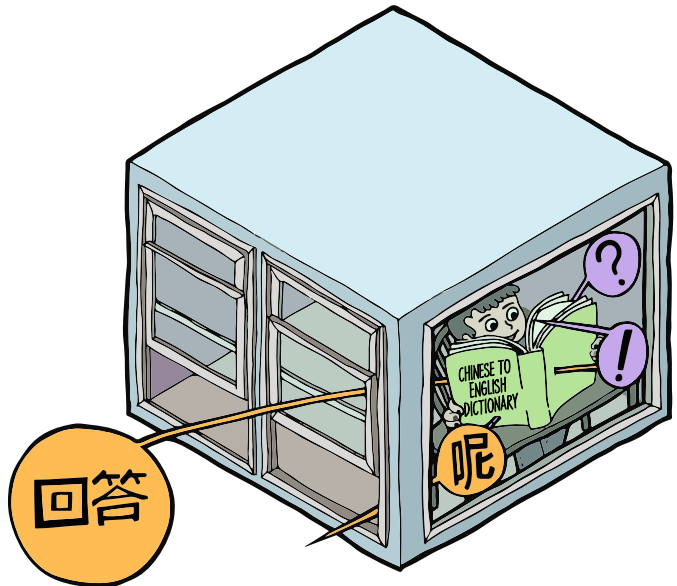
Cover image by Enny Silver and David H. Silver
Cover design by Jeevanjot Kaur Nagpal

**Capish, Com-
prehendes,
Computes?**

Top (Chinese Room

Experiment): The Chinese Room Thought Experiment (John Searle, 1980). An English-speaking person sits inside a sealed room, using a rulebook to manipulate Chinese symbols without understanding their meaning. The system produces fluent-looking Chinese responses purely through syntactic manipulation. Searle's point: following formal symbol-manipulation rules does not constitute understanding.

Bottom (Modern Language Models (LLMs)): Instead of explicit rules, LLMs use high-dimensional vector representations and learned statistical patterns from massive datasets. Inputs are converted into embeddings, processed through multiple nonlinear layers, and decoded into probable outputs. Yet for all the architectural difference, LLMs—like the Chinese Room—operate without intrinsic semantic understanding. They produce contextually appropriate language based on training distributions, not genuine comprehension.



Capish, Comprehendes, Computes?

Searle's Chinese Room thought experiment challenges computational theories of mind: someone manipulates Chinese symbols according to rules without understanding the language. They produce appropriate responses, passing a linguistic Turing test, yet possess no comprehension. The argument distinguishes syntax (symbol manipulation) from semantics (understanding), suggesting that computers executing algorithms operate only at the syntactic level.



CHINESE ROOM EXPERIMENT ◦ SYNTAX WITHOUT
SEMANTICS ◦ SEARLE'S ARGUMENT ◦ STRONG AI
CRITIQUE ◦ LARGE LANGUAGE MODELS ◦ NEXT-WORD
PREDICTION ◦ TRILLION PARAMETERS ◦ COMPRESSION &
GENERALISATION ◦ IN-CONTEXT LEARNING ◦ SUBSTRATE
INDEPENDENCE ◦ BELIEF FORMATION PATTERNS

Words are the only bullets in truth's bandolier.
And poets are the snipers.

— Martin Silenus, circa 2850 A.D.

S'il y avait de telles machines [...] jamais elles ne pourraient user de paroles ni d'autres signes en les composant, comme nous faisons pour déclarer aux autres nos pensées.

("If there were such machines [...] they could never use words by composing them, as we do, to declare their thoughts to others.")

— René Descartes, 1637

Capish, Comprehendes, Computes?

Alan Turing's 1950 essay introduced the 'imitation game' (later known as the Turing Test), proposing that if a machine's textual responses were indistinguishable from a human's, it could be considered intelligent. This marked the beginning of modern debates on machine cognition. In the following decades, the symbolic AI movement gained momentum, with figures such as John McCarthy formalising logic-based systems and Jerry Fodor proposing the 'Language of Thought' hypothesis, which treated mental processes as manipulations of internal symbolic representations.

By the late 1970s, optimism surrounding symbolic approaches began to collide with deeper philosophical questions. Critics questioned whether syntactic manipulation alone could account for semantics—understanding rooted in meaning. In 1980, John Searle articulated the Chinese Room argument, asserting that executing formal rules does not entail comprehension. His critique challenged the assumption of strong AI: that implementing a program is equivalent to having a mind.

Contemporaneously, Daniel Dennett proposed the 'intentional stance,' emphasising observer-relative attributions of belief and intention, while Patricia and Paul Churchland advocated for eliminative materialism, arguing that folk-psychological terms like 'belief' and 'desire' might eventually be replaced by neurobiological accounts. Meanwhile, connectionist models—distributed neural networks—began gaining traction in the mid-1980s, offering an alternative to rule-based systems by emphasising statistical learning over symbolic structure.

By the 1990s, AI had achieved public milestones such as Deep Blue's 1997 victory over Garry Kasparov. Yet critics noted that performance alone does not imply understanding. With the advent of large language models in the 2010s and 2020s, capable of generating coherent and contextually appropriate text, the debate has re-emerged: do these systems understand language, or are they sophisticated instances of the Chinese Room, manipulating symbols without grasping their meaning? Public reports about recent frontier models (e.g., GPT-5) leave specific parameter counts and training token totals undisclosed, though they are trained on massive text corpora at an internet scale.

The Chinese Room thought experiment presents a scenario (Searle, 1980)—an English speaker sits in a sealed chamber, Chinese messages arrive through a slot. The person possesses a rulebook, written in English, that specifies how to manipulate incoming Chinese characters to produce syntactically valid Chinese responses. The rulebook contains no semantic information, only symbol manipulations. The individual follows these instructions and returns the processed strings through the slot.

To an external Chinese speaker, the conversation appears coherent. The responses are grammatically correct, contextually relevant, and indistinguishable from those of a fluent human. Yet the person inside understands none of the content. They do not know that symbols refer to objects, events, or ideas—they execute formal operations on uninterpreted marks. The room satisfies a behavioural test for language competence, yet no part of the system possesses comprehension.

This scenario forces a separation between two dimensions of linguistic behaviour: *syntax*, the arrangement of symbols, and *semantics*, the capacity to represent or grasp meaning. Searle's central claim is that syntactic competence, even when sufficient to pass behavioural tests, does not entail semantic understanding. The system's outputs may simulate language use, but the process lacks intentionality—the directedness of mental states toward meaning-bearing entities or propositions.

This argument extends beyond the thought experiment—it challenges the claims of 'strong AI,' the position that appropriately programmed computers possess minds like humans. Proponents of strong AI maintained that mental states are computational: if a system manipulates symbols according to rules that preserve formal structure and generate appropriate outputs, it qualifies as intelligent. The Chinese Room rejects this inference.

The problem cuts across disciplines. In computer science, the debate centres on algorithmic representation limits and generalisation in machine learning. In linguistics, it intersects with theories of reference, deixis, and semantic grounding. Philosophy confronts questions about intentionality, mental content, and necessary conditions for knowledge. Neuroscience examines embodiment, sensory integration, and causal mechanisms by which mental states arise in biological systems.

From these inquiries emerges a potential requirement: semantic understanding may demand more than pattern matching. Some propose that AI systems might achieve genuine understanding through embodied interaction—robotics, environmental embedding, or sensorimotor coupling that shapes internal representations through causal contact with physical entities. Others argue that meaning resides in subjective experience or first-person perspective that may be inaccessible to artificial systems.

The debate has intensified with large language models (LLMs). These systems demonstrate capabilities that extend far beyond the simple rule-following in Searle's original formulation. They engage in reasoning, exhibit creativity, and show generalisation across domains. Yet they remain neural networks trained through a deceptively simple objective: predicting the next word in a sequence.

The training process operates at an unprecedented scale. Public reports about recent frontier models (e.g., GPT-5) do not disclose exact parameter counts or training token totals; nonetheless, they are trained on massive text corpora at an internet scale. During training, the network processes sequences and learns to predict probability distributions over all possible next words. For the input 'The capital of France is,' the model learns $P(\text{'Paris'}) = 0.85$, $P(\text{'located'}) = 0.03$, and so forth.

This process is self-supervised. No human labels the 'correct' next word because the next word serves as the target. The model minimises cross-entropy loss, heavily penalising confident wrong predictions while providing diminishing returns for improving accurate predictions. Through billions of prediction tasks, spanning months of computation across thousands of processors, the network's parameters converge toward configurations that compress the statistical structure of human language.

At first glance, this training method seems to confirm Searle's critique. The model manipulates symbols based on statistical patterns without direct access to meaning. Critics

dismiss the resulting capabilities as ‘mere probabilistic parroting’: statistical correlation without genuine understanding. This characterisation faces an explanatory challenge that cuts to the heart of the Chinese Room debate. Consider a training example: ‘there are two boxes. Box A contains a red ball and Box B contains a blue ball. If you randomly pick a box and then randomly pick a ball, what is the probability of getting a red ball?’ Now present the model with: ‘there are two containers. Container X holds a cyan sphere and Container Y holds a purple sphere. If you randomly select a container and then randomly choose a sphere, what is the probability of getting a cyan sphere?’

LLMs solve the second problem correctly, yielding 0.5, despite probably never encountering ‘cyan’ and ‘purple’ in mathematical contexts during training. The model abstracts the underlying structure: $P(\text{cyan}) = P(\text{select Container X}) \times P(\text{cyan} \mid \text{Container X}) = 0.5 \times 1.0 = 0.5$. This generalisation cannot be explained by memorization of surface patterns. The specific word combinations or even subsets of this sentence likely never appeared in training data. The model recognises invariant mathematical entities across surface variations, performs conceptual substitution, and transfers zero-shot to novel domains.

This generalisation emerges from a compression constraint. The network must compress terabytes of text into gigabytes of parameters while maintaining prediction accuracy. This compression pressure forces extraction of underlying patterns, rules, and relationships rather than memorization. To predict that ‘The ball rolled down the hill and splashed into the pond,’ the model must develop representations of physics, not just word associations.

Probing techniques confirm this type of learning. Linear classifiers can extract representations of truth, causality, and object properties from the model’s activations. The networks construct hierarchical abstractions. Early layers capture syntax and word boundaries, while later layers encode semantic relationships. This suggests that successful next-word prediction requires building models of the world described in text.

These models undergo multiple training phases that complicate the Chinese Room analogy. Pre-training teaches next-word prediction but produces systems that continue text rather than follow instructions. A model asked ‘What is your first name?’ might respond ‘What is your last name?’, not from understanding but because such sequences appear in training data. Instruction fine-tuning then maps user intentions to appropriate responses through supervised learning on curated instruction-response pairs, transforming text completers into conversational assistants. Finally, reinforcement learning from human feedback drives outputs toward what evaluators judge to be helpful, honest, and harmless.

The specific mechanisms underlying these capabilities matter less than the computational patterns they produce. Current LLMs rely on attention mechanisms, but this appears to be an implementation detail. Alternative architectures achieve similar capabilities through different pathways. This supports substrate independence: intelligence emerges from computational patterns (Putnam, 1967) rather than specific implementations.

This interpretation strengthens when examining capabilities that arise without explicit training. In-context learning allows models to acquire new skills (Brown et al., 2020) from examples in the input prompt, without parameter updates. Present a model with examples of translating English to a made-up language, and it can continue the pattern for new

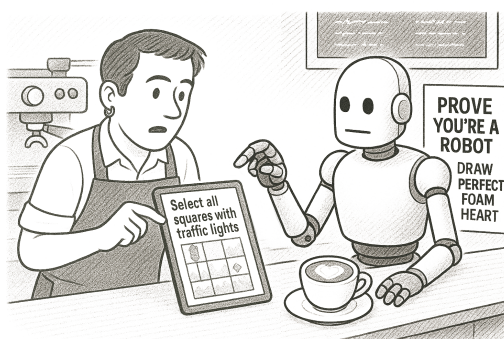
inputs. This suggests that during pre-training, models develop meta-learning algorithms within their forward pass. They maintain implicit probability distributions over possible tasks and update these based on observed examples.

Such capabilities challenge the Chinese Room analogy directly. Searle's scenario involves fixed rule-following. The person executes predetermined instructions without understanding. But LLMs develop adaptive computational patterns that acquire new competencies dynamically. Their 'rules' are not rigid instructions but flexible algorithms that respond to novel contexts. This suggests something different from the static symbol manipulation Searle described.

The core philosophical question persists. The models achieve these capabilities through statistical learning over vast datasets, building representations that compress and generalise from linguistic patterns. Whether this compression constitutes genuine understanding or remains sophisticated simulation is disputed.

Defining 'genuine' versus 'simulated' understanding proves notoriously difficult, yet we possess immediate, non-inferential access to our own comprehension. I know that I grasp meaning—not through behavioural testing or external validation, but through direct phenomenal awareness. Current computational systems lack this quality. The asymmetry is clear. From a first-person perspective, the distinction between understanding and simulation is self-evident. From a third-person perspective, it is non-existent. Other humans and machines occupy the same epistemic position relative to my certainty. I cannot verify understanding in either perspective through observation alone; asserting its existence in other humans requires a leap of faith.

Together with Turing's test (Turing, 1950), the Chinese Room is a test about devising alethic criteria from indistinguishability tests. It forces precision in our definitions of understanding, intelligence, and comprehension. Whether future research reveals that specific architectural features—embodiment, sensorimotor coupling, phenomenal consciousness—are necessary for intelligence, or that substrate independence holds across all cognitive capabilities, Searle's scenario continues to provide a framework for examining these questions.



Cross-Validation.

Belief Formation: Humans and Models

Humans form beliefs in unsettling ways. We think we update our views when presented with new evidence, but reality reveals a different pattern. Beliefs that align with our existing worldview stick around; those that contradict it get dismissed or twisted into supporting evidence. When someone challenges our deep convictions, we often become *more* confident in what we believed originally—sometimes described as the **backfire effect** (though evidence suggests such effects are not widespread and are context-dependent). We're not neutral fact-processing machines. We're defensive storytellers, preserving narrative coherence over empirical accuracy.

We can now compare this to large language models. When you train a model on billions of text examples, it learns whatever patterns exist in that data, including confident assertions about false claims. Later, when researchers try to fine-tune the model with correct information, the original learning resists change. The neural weights have settled into configurations optimised for the original data distribution.

This is how learning systems work. Both human brains and neural networks must compress vast amounts of information into manageable models. Once those patterns solidify, changing them means destabilising everything else that depends on them. In humans, this shows up as cognitive dissonance and motivated reasoning. In models, it appears as **gradient stasis** and **catastrophic forgetting**. This shows up as cognitive gradient stasis when learning new ones.

Both systems handle uncertainty similarly. Humans rarely admit ignorance cleanly. Instead, we confabulate. Language models do the same. When asked about topics outside their training data, they don't respond with "I don't know." They generate confident-sounding responses that preserve conversational flow, even when information is sparse or contradictory.

This suggests that both human cognition and current AI systems are optimised for something other than truth correspondence. They prioritise internal consistency and social coordination over factual accuracy. In humans, this makes evolutionary sense. Being wrong together was often more adaptive than being right alone. In AI systems, it emerges from the training objective: predict the next word in a way that sounds human-like.

The Chinese Room becomes more provocative through this lens. Searle asked whether symbol manipulation without understanding constitutes genuine comprehension. But perhaps the more unnerving question is whether human 'understanding' is itself merely symbol filtering that prioritises narrative coherence over external reality.

The Epistemic Fragility of Syntax-Only Cognition

Framing the Dispute

The Chinese Room is not an empirical argument but a methodological critique aimed at conceptual boundaries in cognitive science. Searle's challenge is not about empirical performance but what we are permitted to infer from it. He contends that passing the Turing Test—or any test based solely on linguistic output—cannot entail genuine understanding unless we already have a theory that licences such an inference. The argument's strength lies in its methodological reversal: where AI seeks to move from behavioural evidence to mental attribution, Searle denies the validity of that inference without a prior grounding in what it means to 'understand'.

Understanding Without a Criterion

'Understanding' is not an operationally neutral term. Unlike 'predicts weather' or 'stores data,' it is irreducibly normative and semantically loaded. To assert that a system understands is to claim it possesses internal relation to meaning—content-bearing capacity that cannot be exhausted by structural descriptions alone. The Chinese Room exposes justificatory laziness: we project understanding onto systems that behave in familiar ways, without specifying what justifies that projection.

Replies to Searle rely on stipulative bridges—identifying understanding with functional role, environmental embedding, or counterfactual dependence. These bridges merely relocate the conceptual burden. What is lacking is a non-circular account of when formal behaviour amounts to semantic content. This is not a technical failure but a philosophical silence.

Intentionality and Attribution

Searle's intentionality point is often misconstrued. He is not claiming that syntax cannot, in principle, be paired with semantics. He asserts that such pairing is not guaranteed by formal operations alone. The Chinese Room Argument (CRA) is not about what symbols do; it is about what they mean. And meaning

is not intrinsic to the system unless some internal state stands in a relation of intentional directedness—a relation not captured by computational transitions.

Attempts to circumvent this by pointing to system-wide properties (as in the Systems Reply) or virtual entities (as in the Virtual Mind Reply) fail to address a fundamental asymmetry: intentional states have first-person authority, whereas syntactic states do not. If a system 'understands,' then it makes sense to ask what it understands and why. But if that attribution is based only on labelling (e.g., 'this system understands Chinese'), we are no longer explaining cognition—merely renaming behaviour.

Simulation and Normativity

Searle targets the normative dimension of cognition. Understanding involves norm-sensitive responsiveness to content, not merely causal states. A person can misunderstand, misinterpret, or revise their understanding. These are not errors in computation; they are errors in relation to content. But a syntactic machine cannot err in this sense. It can malfunction, but it cannot misbelieve. Without normativity, there is no epistemic traction, and without that, no understanding.

The Epistemic Cost of Ambiguity

The enduring appeal of the Chinese Room stems from its methodological clarity. It does not claim that AI will never understand. It claims that we lack a criterion by which to know if it does. To assert that future systems might understand language is to talk without terms. Until we have a definition of understanding that does not collapse into performance, or a theory of meaning that does not presume biological embedding, our attributions remain projections—not findings.

References:

- Searle, J. R. (1980). *Minds, Brains, and Programs*. Behavioral and Brain Sciences, 3(3), 417–457.
- Dennett, D. C. (1987). *The Intentional Stance*. MIT Press.
- Chalmers, D. (1996). *The Conscious Mind*. Oxford University Press. See Also: <https://plato.stanford.edu/entries/chinese-room/>

