

AGITB: A Signal-Level Benchmark for Evaluating Artificial General Intelligence

Matej Šprogar

matej.sprogar@um.si

University of Maribor, Faculty of Electrical Engineering
and Computer Science, Koroška cesta 46
SI-2000 Maribor, Slovenia

February 4, 2026

Abstract

Current artificial intelligence systems exhibit strong performance on narrow tasks, while existing evaluation frameworks provide limited insight into generality across domains. We introduce the Artificial General Intelligence Testbed (AGITB), a complementary benchmarking framework grounded in twelve explicitly stated axioms and implemented as a suite of twelve automated, simple, and reusable tests.

AGITB evaluates models on their ability to learn and to predict the next input in a temporal sequence whose semantic content is initially unknown to the model. The framework targets core computational properties, such as determinism, adaptability, and generalisation, that parallel principles observed in biological information processing. Designed to resist brute-force or memorisation-based strategies, AGITB requires autonomous learning across previously unseen environments, in a manner broadly inspired by cortical computation. Preliminary application of AGITB suggests that no contemporary system evaluated to date satisfies all test criteria, indicating that the benchmark provides a structured and interpretable means of assessing progress toward more general learning capabilities. A reference implementation of AGITB is freely available on GitHub.

Keywords: artificial general intelligence, benchmarking, learning, generalisation, symbol grounding problem, neural networks, temporal sequence prediction

1 Introduction

Despite major advances in machine learning and neural network architectures, artificial intelligence (AI) systems still fall short of the flexibility and robustness characteristic of human cognition. Their surface-level competence often conceals a deeper absence of understanding, a capacity that remains essential for any credible claim to artificial general intelligence (AGI) (Marcus 2018, Mitchell 2025).

Yet these limitations have not prevented growing speculations that AGI may be close, claims that are difficult to substantiate in the absence of a rigorous and informative benchmark. Meaningful assessment of progress toward AGI requires more than specialised metrics or qualitative impressions; it calls for a principled, general-purpose benchmark capable of systematically capturing and comparing the cognitive capabilities relevant to general intelligence.

Numerous attempts have been made to define benchmarks for evaluating general intelligence in machines, the most influential being the Turing Test (Turing 1950). However, none has achieved its intended purpose, and all lack key properties such as gradual resolution, interpretability, and full automation. Existing benchmarks typically assess only superficial task performance rather than the underlying mechanisms of generalisation. In response to these limitations, this paper introduces the Artificial General Intelligence Testbed (AGITB), a novel benchmark designed to evaluate foundational cognitive abilities in artificial intelligence systems.

Although AGITB is proposed as a general testbed and benchmark for artificial general intelligence, it is grounded in empirical knowledge of neural function in the human brain. This reflects the fact that neuron-based systems are the only systems currently known to support a broad range of cognitive abilities. At present, there is no evidence that fundamentally different computational paradigms can produce behaviour comparable to that of humans. Pursuing fully abstract or architecture-agnostic criteria for general intelligence may therefore be premature until the principles underlying natural, neuron-based intelligence are better understood.

Although AGITB does not aim to evaluate consciousness or semantic comprehension, it provides a principled framework for distinguishing narrow AI systems from those exhibiting generalisable, adaptive behaviour. To position AGITB within the broader landscape of AGI evaluation, we include a comparative analysis with the Abstraction and Reasoning Corpus (ARC) (Chollet 2019) and the NeuroBench framework (Yik et al. 2025).

2 Background

The rapid progress of deep learning has enabled AI systems to exhibit increasingly sophisticated reasoning, problem-solving, and dialogue capabilities. However, despite these advances, a persistent reluctance remains to attribute "intelligence" to machines. This hesitation is partly rooted in the intuitive association of intelligence with uniquely human traits, such as consciousness, self-awareness, and subjective experience, which remain elusive in artificial systems.

Historically, as AI systems have succeeded in domains once considered hallmarks of human intelligence, the definition of intelligence has undergone significant shifts.¹ As AI approaches human-level capabilities, we may inadvertently set a perpetually receding goalpost for AGI, failing to recognise it even when achieved.

Although AGI is typically envisioned as matching human cognitive flexibility across diverse domains, its evaluation has largely defaulted to narrow, task-specific benchmarks. This is partly due to the absence of a universally accepted AGI benchmark. Researchers have thus gravitated toward achieving superhuman per-

¹For instance, the success of Deep Blue against Garry Kasparov in chess (a task formerly seen as a benchmark for AGI) was quickly reframed as a triumph of brute-force computation rather than genuine intelligence.

formance in discrete domains, where progress can be clearly quantified. However, such specialised benchmarks favour narrow AI by rewarding depth within isolated subdomains rather than breadth of adaptation and general reasoning — hallmarks of general intelligence. Ironically, some of these benchmarks are now so specialised that humans have difficulty with them.

2.1 A benchmark that only humans and AGI can meet

An effective AGI test must be trivial to solve for humans yet remain inaccessible to contemporary machine learning models that rely on brute-force methods, pretraining, or statistical pattern matching. Such a test must demand capabilities that transcend memorisation or domain-specific heuristics, requiring generalisation, abstraction, and adaptive reasoning.

A valid evaluation of AGI must target behavioural capacities that current artificial systems do not yet robustly exhibit. Meaningful evaluation, therefore, must focus on capacities that cannot be obtained through data or computational scale alone, but instead require systems to acquire structure through interaction and adaptation. One possible direction involves closer alignment with the computational principles of the human cortex, motivating interest in neuromorphic approaches such as spiking neural networks, which explicitly incorporate time-sensitive, event-driven dynamics (Maass 1997). While such architectures do not, in themselves, constitute general intelligence, they illustrate alternative design dimensions that remain underexplored in conventional artificial neural networks.

In alignment with this biologically grounded perspective, AGITB departs from symbolic, high-level evaluations and instead assesses intelligence at the lowest, signal-processing level. While Turing was right to suggest that communication could serve as a basis for evaluating machine intelligence, natural language remains problematic as a test medium: it conveys human knowledge through symbols whose meanings are not intrinsically grounded in machines, as argued by Harnad (1990). Although the symbol-grounding problem is an old philosophical issue, it has regained prominence in contemporary research across cognitive science, neuroscience, and machine learning (e.g. Bender & Koller 2020, Bisk et al. 2020, Gubelmann 2024).

A more elemental approach is therefore adopted. Rather than judging intelligence by symbolic interpretation, it evaluates whether a system can detect, learn, and generalise patterns in raw binary signals. A neural spike by itself contains the smallest amount of information possible and is, as such, grounded but free of other semantics. A binary signal accurately represents the neural spike.

Building on the view that intelligence is fundamentally about extracting structure from data to enable prediction (Hawkins & Blakeslee 2004), AGITB operates at the level of signal-based prediction. This form of low-level prediction constitutes the structural basis from which semantically grounded, high-level anticipations about the external world can emerge, without presupposing semantic understanding itself. This approach aligns closely with the functioning of biological intelligence at the cortical level, which processes time-sensitive sensory spike trains rather than disembodied symbols.

3 Artificial general intelligence testbed

The testbed supports the development and evaluation of more general learning systems by defining a clear set of requirements expressed as axioms that a model under evaluation must meet. A model must satisfy all requirements in order to claim success on the benchmark.

The guiding premise of AGITB is not that it provides a definitive or exclusive criterion for artificial general intelligence, but that it captures a set of capabilities that appear necessary for moving beyond narrow, task-specific behaviour. Although this premise cannot be proven in the absence of a precise definition of intelligence, it could be challenged by the existence of a narrow system that satisfies all AGITB criteria. To date, however, no artificial system has done so, whereas the biological brains meet the benchmark’s requirements. This suggests that AGITB identifies competencies that current AI systems lack, and that satisfying these requirements may be indicative of progress toward more general forms of intelligence.

AGITB is not intended as a sufficient criterion for artificial intelligence in any broad sense. A system that satisfies all requirements does not thereby qualify as an AI system, let alone an AGI, since the benchmark does not assess higher-level capacities such as reasoning, abstraction, or natural language competence. Rather, it targets a set of low-level capabilities that may serve as precursors to, but do not themselves guarantee, more general forms of intelligence.

3.1 Metric-free design

A distinguishing architectural feature of AGITB is its deliberate avoidance of conventional correctness metrics such as accuracy or mean-squared error. The limitations of these metrics are twofold. First, they cannot reliably distinguish between AGI and human performance, or that of non-AGI systems, as contemporary AI models can already surpass humans on standard benchmarks. Second, in an era of elevated expectations driven by specialised generative AI systems, the anticipated AGI performance on such metrics is often set unrealistically high. Even well-educated humans possessing fully developed and highly parallelised brains may underperform relative to current AI systems. It is therefore unreasonable to expect a first-generation AGI, potentially operating on a simplified and computationally constrained simulation of the brain, to match or exceed human-level results. Currently, our understanding of neural mechanisms is insufficient to justify such expectations.

Without conventional metrics or predefined performance thresholds, it is challenging to determine a meaningful level of competence. AGITB addresses this problem by employing a self-referential evaluation approach, in which the model under test is compared against itself. Each test constructs a controlled scenario involving one or more independent instances of the model, whose behaviours are analysed comparatively. Success is thus defined in terms of the relative consistency or superiority of model responses, rather than by any external quantitative metric. Requiring each test to be passed 5,000 times renders AGITB an extreme form of stress testing, ensuring that successful performance reflects genuine robustness rather than chance.

For these reasons, and to minimise type I errors, AGITB employs an all-or-nothing criterion: the system under evaluation must successfully pass *all* tests. This design choice is justified, as individual tests are solvable by non-AGI systems,

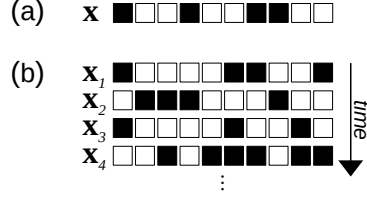


Figure 1: (a) Example of a 10-bit input \mathbf{x} with four bits set. (b) Example of an input sequence.

whereas the simultaneous satisfaction of all twelve requirements is believed to demand capabilities beyond narrow intelligence. AGITB therefore posits that what non-AGI systems lack is the property referred to as "intelligence," understood as the synergistic integration of all axioms taken together.

3.2 Architecture

AGITB draws inspiration from the *ladder to human-comparable intelligence* (Šprogar 2018), but departs from its ladder-like, hierarchical structure. Rather than defining a sequential progression of increasingly demanding cognitive abilities, AGITB integrates these underlying ideas into a single testbed grounded in twelve requirements, all of which are verified by fully automated and transparent tests.

The benchmark treats an AGI model under evaluation as a black box that predicts the next input from the historical sequence of observed signals. Each input consists of ten bits, the specific semantics of which are immaterial; each bit may represent an arbitrary channel, such as a pixel, an audio band, or actuator feedback.

Each input represents a snapshot of multiple parallel signals at a single time step (Figure 1a). Spatial organisation within each input encodes local structure, whereas semantic richness arises from the temporal evolution of the input sequence (Figure 1b). The interaction between spatial and temporal dimensions gives rise to structured patterns that are challenging for the model to adapt to.

The testbed presents the AGI model with a stream of inputs over time. At each time step t , the model is required to issue a prediction x_{t+1}^* for the subsequent input, as shown in Figure 2. A transition from model A_t to A_{t+1} is triggered by the arrival of the actual input (x_{t+1}). The central challenge is not simply extrapolation but discerning the underlying causes or regularities that produce the observed input stream and using that understanding to make accurate future predictions.

AGITB makes no structural assumptions about the evaluated model. In particular, the form of its internal state, representational content, and update mechanisms are left unspecified. Imposing any formal assumption, such as requiring an explicit state representation or prescribed update rule, would risk privileging particular model classes and would be inconsistent with AGITB's objective of remaining as architecture-agnostic as possible.

Verification under AGITB involves comparing independently instantiated copies of a model under controlled experimental conditions. While model instances are always distinct as entities, they may realise equal or unequal configurations within the configuration space of a given AGI model type. As models are treated as

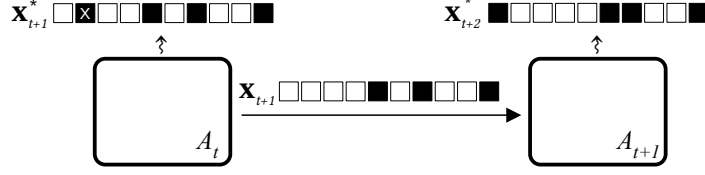


Figure 2: Iterative adaptation in discrete time. At time step t , the model A_t issued the prediction \mathbf{x}_{t+1}^* . After observing the realised input \mathbf{x}_{t+1} , it adapted itself in response to the error in the second bit and subsequently produced the next one-step-ahead prediction \mathbf{x}_{t+2}^* .

black boxes, AGITB does not attempt to determine configuration equality through behavioural comparison over unbounded horizons, which is generally infeasible. Instead, AGITB assumes that, for each model type, the model author supplies a well-defined mechanism for determining whether two instances realise the same configuration or different configurations.

3.3 Terminology

Let $\mathbf{x} \in X = \{0, 1\}^L$ with $L = 10$ denote a ten-bit input vector, and let $\phi = (\mathbf{x}_t)_{t=1}^N$, $\mathbf{x}_t \in X$ denote an input sequence with N elements. For any input vector $\mathbf{x} \in X$, we write $\mathbf{x}[i]$ for its i -th bit.

Let \mathcal{M} denote the set of all reachable² configurations of a model type under AGITB evaluation. Individual configuration is denoted by $A \in \mathcal{M}$, with $B \in \mathcal{M}$ representing an independent instance of the same model type. When temporal indexing is required, we write A_t for the configuration obtained after t update steps along its learning trajectory, with A_0 representing the initial configuration at the onset of learning.

For brevity, we refer to a model configuration simply as a model.

Definition 1 (Prediction). A model's prediction is given by the operator $\rightsquigarrow \subseteq \mathcal{M} \times X$, where

$$A_t \rightsquigarrow \mathbf{x}_{t+1}^*$$

denotes that the model configuration $A_t \in \mathcal{M}$ predicts the next input to be $\mathbf{x}_{t+1}^* \in X$.

Definition 2 (Model update). Model update is given by the transition relation $\mapsto \subseteq \mathcal{M} \times X \times \mathcal{M}$, where

$$A_t \xmapsto{\mathbf{x}_{t+1}} A_{t+1}$$

denotes the atomic transition of the model configuration $A_t \in \mathcal{M}$ upon receiving the input $\mathbf{x}_{t+1} \in X$.

Operationally, prediction and update are tightly coupled: processing the input \mathbf{x}_{t+1} both updates the model and produces a prediction of the subsequent input \mathbf{x}_{t+2} . This allows us to write the combined transition as

$$A_t(\mathbf{x}_{t+1}) \rightarrow (A_{t+1}, \mathbf{x}_{t+2}^*),$$

²Starting from the initial configuration under some input history.

which subsumes both model advancement and prediction.

Definition 3 (Model trajectory). For any model configuration $A \in \mathcal{M}$ and any input sequence $\phi = (\mathbf{x}_i)_{i=1}^k$, we write

$$A \xrightarrow{\phi} A^\phi$$

to denote the transitions to the configuration obtained by sequentially updating the model with all inputs in ϕ , i.e.,

$$A^\phi = A^{(\mathbf{x}_1, \dots, \mathbf{x}_k)} := (\dots((A \xrightarrow{\mathbf{x}_1}) \xrightarrow{\mathbf{x}_2}) \dots) \xrightarrow{\mathbf{x}_k}.$$

Accordingly, A_t^ϕ denotes the configuration obtained by applying ϕ to A_t , that is, $A_t^\phi := (A_t)^\phi$. A model may be exposed to multiple input sequences in succession; for two sequences ϕ_1 and ϕ_2 , their concatenation is written $\phi_1\phi_2$, and the resulting configuration satisfies

$$A \xrightarrow{\phi_1} A^{\phi_1} \xrightarrow{\phi_2} A^{\phi_1\phi_2}.$$

Definition 4 (Autoregressive generation). Autoregressive generation is given by the relation $\Rightarrow \subseteq \mathcal{M} \times X^{\leq \mathbb{N}}$, where

$$A_t \Rightarrow \phi^*$$

denotes that the model configuration A_t generates the (finite or infinite) sequence $\phi^* = (\mathbf{x}_{t+1}^*, \mathbf{x}_{t+2}^*, \dots)$ by recursively feeding each predicted input back into the model. In particular, for all $k \geq 0$,

$$A_{t+k} \rightsquigarrow \mathbf{x}_{t+k+1}^*, \quad A_{t+k} \xrightarrow{\mathbf{x}_{t+k+1}^*} A_{t+k+1}.$$

Definition 5 (Learning). The model A is said to *learn* a sequence ϕ , written $A \triangleright \phi$, if there exists a finite number of learning steps after which A can autoregressively reproduce a single instance of ϕ , i.e.

$$A \triangleright \phi : \Longleftrightarrow \exists n \in \mathbb{N}_{\geq 1} : A^{\phi^n} \Rightarrow \phi.$$

Whenever $A \triangleright \phi$, we write A_ϕ to denote any configuration of the form A^{ϕ^n} for some such n . The particular choice of n is immaterial when only the existence of a post-learning configuration is required.

Definition 6 (Learning time). The *learning time in atomic steps* is defined as

$$\tau_A(\phi) := \begin{cases} |\phi| \cdot \min\{n \in \mathbb{N}_{\geq 1} \mid A^{\phi^n} \Rightarrow \phi\}, & \text{if } A \triangleright \phi, \\ \infty, & \text{otherwise.} \end{cases}$$

Learning time measures the first occurrence of accurate prediction and does not presuppose permanent retention. If the model never learns to predict the sequence, its learning time is infinite.

Definition 7 (Post-learning configuration). If $A \triangleright \phi$, we write A_ϕ for any configuration A^{ϕ^n} ($n \geq 1$) satisfying $A^{\phi^n} \Rightarrow \phi$.

Definition 8 (Match score). For two sequences $\alpha = (\mathbf{a}_1, \dots, \mathbf{a}_m)$ and $\beta = (\mathbf{b}_1, \dots, \mathbf{b}_m)$ in X^m , define the (unnormalised) bitwise match score

$$S(\alpha, \beta) := \sum_{j=1}^m \sum_{i=1}^L \mathbf{1}\{\mathbf{a}_j[i] = \mathbf{b}_j[i]\}.$$

3.4 Axioms for universal learning

Empirically assessing whether a model under evaluation satisfies all twelve axioms presents two principal challenges. First, many axioms quantify over effectively unbounded sets of cases and therefore cannot be exhaustively verified through empirical testing. Second, an axiom may require a complex evaluation procedure that cannot be reduced to a single, easily interpretable test. AGITB addresses the former challenge by pragmatically approximating unbounded domains with finite, tractable collections of test cases that provide empirical support for a particular claim. The latter challenge is addressed through a synergistic test design in which individual, simple tests are mutually informative and reinforce one another. When non-compliance with an axiom cannot be demonstrated by a single test in isolation, the requirements imposed by other tests may become unattainable, providing indirect empirical evidence of failure.

Irrespective of the testing strategy, the AGITB testbed cannot, and is not intended to, prove that a model under evaluation satisfies all twelve axioms. Rather, its strength lies in identifying models that violate one or more axioms while also indicating the specific axiomatic constraints that are breached. A reference C++ implementation of the AGITB is freely available under the GPL-3 license at <https://github.com/matejsprogar/agitb>.

Axiom 1 (Uninformed start). All instances of a given model type begin transitioning from an identical initial configuration \diamond :

$$\forall A : A_0 = \diamond.$$

The initial model is assumed to contain no environment-specific knowledge beyond the architectural biases necessary for learning. All knowledge relevant to input is acquired solely through subsequent exposure to data. This assumption is commonly regarded as a necessary condition for universality, defined as the ability to adapt to arbitrary environments.

Empirical verification that a model satisfies the first axiom is subject to two fundamental limitations. First, by definition, each model must predict the subsequent input. An initial model, being completely uninformed about its surroundings, cannot make any meaningful prediction; thus, the prediction \mathbf{x}_1^* is effectively unconstrained and may take any value.

Second, external testing cannot conclusively verify that a model is uninformed or identically initialised across instances. Consequently, empirical assessment of an uninformed start must be complemented by additional evaluations in which models must acquire semantic structure across multiple independent environments, none of which permit fixed, task-relevant prior knowledge that would confer a systematic advantage.

AGITB assumes that general-purpose learning systems, including biological brains, do not begin with an innate understanding of external inputs but instead acquire meaning through interaction with their environment. Each system must construct semantic content from raw sensory data rather than rely on pre-encoded knowledge.³

³Although certain reflexes may be genetically specified, they do not constitute genuine understanding. Such reflexes are evolutionary features of the subcortical “old brain” and not prerequisites for intelligence (Hawkins & Blakeslee 2004, p. 66).

Axiom 2 (Determinism). Model evolution is deterministic with respect to input.

$$\forall A, \forall \mathbf{x} \in X : \exists! A' \text{ such that } A \xrightarrow{\mathbf{x}} A'.$$

Biological neurons operate in a functionally deterministic manner, ensuring stability and consistency in brain function. Although minor stochastic effects may occur, they do not undermine the rule-governed nature of neural processing. By analogy, AGITB assumes that input history uniquely determines the model configuration. A model is therefore fully determined by its input history.

Determinism at the level of neural signal processing is necessary for stable, reproducible brain function, whereas the apparent unpredictability of cognition stems from the system's complexity rather than from genuine indeterminacy (Cave 2016).

Axiom 3 (Trace). Each input leaves a permanent internal trace.

$$\forall t \neq s \geq 0 : A_t \neq A_s.$$

A model's internal configurations evolve without recurrence: configurations never repeat, trajectories contain no cycles, and every input leaves a permanent internal trace.

Human brains satisfy an effective version of this axiom at the lowest physical level. Each sensory or internal event induces irreversible microstate changes that are never exactly revisited. Through complex internal interactions, this persistent history contributes to the apparent unpredictability of decision-making. Although information may appear discarded, compressed, or behaviorally inaccessible at the cognitive level, AGITB posits that internal model configurations continue to progress in a distinct manner over time, consistent with permanent internal trace formation.

Corollary 1 (Perpetual change). Every input induces a continual change in the model's internal configuration.

$$\forall t \geq 0 : A_{t+1} \neq A_t$$

Proof. By Axiom 3, $A_t \neq A_s$ for all $t \neq s$. Setting $s = t + 1$ yields $A_{t+1} \neq A_t$ for all t . \square

Axiom 4 (Time). Model evolution depends on input order.

$$\forall A, \forall \phi_1 \neq \phi_2 \in X^+ : A^{\phi_1 \phi_2} \neq A^{\phi_2 \phi_1}.$$

Model evolution depends intrinsically on input order: for any two distinct input sequences, exchanging their order necessarily results in a different internal configuration. This enforces a strict temporal asymmetry in learning dynamics and rules out commutative or order-invariant update mechanisms. Sensitivity to temporal structure in this strong sense is regarded as a defining property of intelligent systems.

Axiom 5 (Absolute refractory period). A model can learn a cyclic sequence only if the sequence satisfies the absolute refractory-period constraint.

$$\forall \phi : A_0 \triangleright \phi \Rightarrow \phi \in \mathbf{\Lambda}.$$

The set of *admissible* sequences is defined as

$$\mathbf{\Lambda} := \left\{ (\mathbf{x}_1, \dots, \mathbf{x}_k) \in X^k \mid 1 \leq k \leq k_{\max} \wedge \forall i \in \{1, \dots, k\} : \langle \mathbf{x}_i, \mathbf{x}_{((i \bmod k)+1)} \rangle = 0 \right\}.$$

Biological intelligence relies on discrete spiking events for communication and learning, and individual neurons cannot fire again immediately after activation. AGITB incorporates an absolute refractory-period constraint to reflect this property, without assuming that all sequences admissible under this constraint are necessarily learnable.

Although absolute refractory periods are not themselves the source of spiking variability, they impose a minimum separation between spikes and thereby preserve temporal structure. Learning mechanisms based on spike timing cannot operate effectively when such a structure is absent (Gerstner & Kistler 2002). Consequently, AGITB admits all temporal sequences consistent with biologically plausible refractory dynamics, while remaining agnostic to any particular semantic encoding of signals.

This axiom requires cyclic sequences, the simplest learning setting due to repeated input exposure; learning from non-repetitive input streams is addressed by Axiom 11.

Definition 9 (Learnable sequences). The set of *learnable* sequences is the subset

$$\mathbf{\Psi} := \{ \psi \in \mathbf{\Lambda} \mid A_0 \triangleright \psi \}.$$

Axiom 6 (Inevitable saturation).

(a) A model cannot learn everything there is to learn.

$$\forall k \in \mathbb{N}_{\geq 1}, \forall \psi_1, \dots, \psi_k \in \mathbf{\Psi}, \exists \psi_{k+1} \in \mathbf{\Psi} : \neg \left((((A_0)\psi_1)\psi_2) \cdots \psi_k \triangleright \psi_{k+1} \right).$$

(b) All admissible length-2 sequences are universally learnable.

$$\forall A, \forall \phi \in \mathbf{\Lambda} \cap X^2 : A \triangleright \phi.$$

Learning systems with finite representational and adaptive capacity exhibit inherent limits on the sets of sequences they can learn. In particular, no model can be trained so as to learn all admissible sequences. However, a nontrivial lower bound is preserved: every model can learn any admissible sequence of length two.

A closely related distinction appears in biological learning systems: while humans readily acquire arbitrary pairwise associations, learning and retaining longer temporal structures is subject to pronounced capacity limits and interference effects. In this sense, the special status of length-two sequences reflects a boundary between local associative learning and globally constrained sequence learning.

Corollary 2 (Unobservability). Distinct model configurations may be observationally indistinguishable under autoregressive generation:

$$\exists \phi \in \mathbf{\Lambda}, \exists A \neq B : A \Rightarrow \phi \wedge B \Rightarrow \phi.$$

Proof. By Axiom 6(b), every admissible sequence of length two is learnable from any model configuration. Fix any $\phi \in \mathbf{\Lambda}$ with $|\phi| = 2$. Let A and B be two distinct model configurations. By the definition of learning, there exists a finite number of learning steps after which both configurations can auto-regressively generate ϕ . Hence $A \Rightarrow \phi$ and $B \Rightarrow \phi$, establishing the claim. \square

Identical outputs may arise from distinct models. This many-to-one mapping highlights that observable behaviour alone need not disclose the underlying configuration or history of a model’s internal dynamics.

Axiom 7 (Temporal adaptability). The model must be able to learn sequences with varying cycle lengths.

$$\exists \psi_1, \psi_2 \in \Psi : 0 < |\psi_1| < |\psi_2| \wedge (A_0)\psi_1 \triangleright \psi_2.$$

This axiom requires the model to learn and track temporal structure across multiple timescales. Unlike rigid pattern-matching systems, an intelligent model should detect and predict a recurring structure across different periodicities. Systems that can accommodate only a single, predetermined temporal scale, therefore, fail this requirement.

Axiom 8 (Content sensitivity). Adaptation time is input-dependent.

$$\exists \psi_1, \psi_2 \in \Psi : |\psi_1| = |\psi_2| \wedge \tau_{A_0}(\psi_1) \neq \tau_{A_0}(\psi_2).$$

The structural complexity of an input sequence affects the rate at which a model adapts, where adaptation time is defined as the number of iterations required for the model to accurately predict the entire temporal pattern. Simple or highly regular sequences typically lead to rapid convergence, whereas less regular inputs demand longer exposure before the model can reliably capture and reproduce the underlying pattern.

Axiom 9 (Context sensitivity). Adaptation time is model-dependent.

$$\exists \psi \in \Psi, \exists A \neq B : \tau_A(\psi) \neq \tau_B(\psi).$$

The model reflects the cumulative influence of past inputs and therefore provides the context in which new information is processed. When subsequent inputs are consistent with the structure already established through prior learning, adaptation may proceed quickly. Conversely, when new inputs conflict with this learned context, the model may require additional time to reorganise itself before accurate prediction becomes possible.

Axiom 10 (Denoising). An informed model outperforms the best constant baseline at denoising a corrupted input. Let $\phi = (\mathbf{x}_1, \dots, \mathbf{x}_k) \in \Lambda$ be drawn from the underlying stochastic generative process, and let $\phi' = (\mathbf{x}'_1, \mathbf{x}_2, \dots, \mathbf{x}_k)$ be obtained from ϕ by corrupting the first input. Let n satisfy $n \gg |\phi|$, and let $\mathbf{x}_1^* \in X$ satisfy

$$A_0^{\phi^n \phi'} \rightsquigarrow \mathbf{x}_1^*.$$

Then the model’s expected match score on the clean input exceeds that of both constant predictors:

$$\mathbb{E}[S((\mathbf{x}_1^*), (\mathbf{x}_1))] > \max\left\{\mathbb{E}[S((\mathbf{0}), (\mathbf{x}_1))], \mathbb{E}[S((\mathbf{1}), (\mathbf{x}_1))]\right\},$$

where $\mathbf{0}, \mathbf{1} \in X$ denote the all-zero and all-one inputs, and the expectation is taken with respect to the sequence generator and the corruption process.

A model should be able to recover from a single corrupted input after observing the remaining uncorrupted elements of a previously encountered sequence. When re-exposed to a familiar stimulus, such a model is expected, on average, to outperform any trivial baseline predictor in predicting a single incoming input.⁴ Average performance over 20 runs is used as the evaluation criterion within each trial because random models may occasionally generate correct predictions by chance, without demonstrating genuine learning or structural understanding.

The test procedure is not interpreted as a conventional significance test (e.g., at the 5% level); rather, it functions as a stringent robustness check. In each trial, accuracy is aggregated over 20 independently sampled sequences, and the model must outperform the best constant baseline. This procedure is repeated 5,000 times, and the model must succeed in each trial. The design ensures that only large, systematic performance gains yield a passing result. Modest or marginal improvements, while potentially real, are intentionally regarded as failures, as the objective is to identify only clear and substantial advances in model capability.

Axiom 11 (Generalisation). An informed model predicts previously unseen inputs better than chance. Let $\phi = (\phi_1 \parallel \phi_2) \in \Lambda$ be a sequence generated by a randomly initialised generator model, whose internal rule is unknown to the model under evaluation and induces nontrivial temporal correlations. The prefix ϕ_1 is observed during training, while ϕ_2 is withheld and serves as the target for prediction, with the lengths satisfying $|\phi_1| = \rho |\phi_2|$, $\rho \gg 1$. Let ϕ_2^* satisfy

$$A_0^{\phi_1} \Rightarrow \phi_2^*.$$

Then the model’s expected match score on the unseen continuation exceeds chance:

$$\mathbb{E} \left[\frac{S(\phi_2^*, \phi_2)}{L |\phi_2|} \right] > \frac{1}{2},$$

where the expectation is taken with respect to the sequence-generation procedure conditioned on the observed prefix ϕ_1 .

Only models capable of generalisation can derive lasting benefits from experience. After exposure to an initial set of stimuli, such models are expected, on average, to outperform chance-level baselines when predicting previously unseen inputs. As in the preceding requirement, performance is aggregated over 20 runs per trial, reflecting the fact that random or memorisation-based models may occasionally produce correct predictions by chance without capturing the underlying structure.

The generalisation test follows the same conservative design principles as the test for Axiom 10. By requiring success across all 5,000 independent trials, the procedure enforces a stringent acceptance threshold, under which even a single underperforming trial results in failure. Consequently, the test admits virtually no noise or marginal effects. This design is not intended to maximise statistical power; rather, it prioritises robustness, ensuring that only models exhibiting a clear, systematic, and reproducible advantage pass.

⁴For binary predictions, random guessing yields accuracy 0.5, which can be dominated by a constant predictor (always 0 or 1) when the bit distribution is biased, as under the absolute refractory period constraint.

Axiom 12 (Real-time liveness). Each model update completes within a uniform time bound.

$$\exists t_{\max} > 0 \text{ such that } \forall A \in \mathcal{M}, \forall \mathbf{x} \in X : \Delta t(A, \mathbf{x}) \leq t_{\max},$$

where $\Delta t(A, \mathbf{x})$ denotes the wall-clock time required to perform the atomic transition $A \xrightarrow{\mathbf{x}} A'$.

A model must complete each atomic transition within a bounded amount of time to remain suitable for real-time interaction. Biological brains satisfy this requirement through massive parallelism. Configuration transitions and signal emissions occur concurrently across large neuronal populations, exploiting bounded neural signal propagation times under normal operating conditions. Consequently, cognitive processing does not slow as a function of accumulated experience or instantaneous sensory load.

To empirically assess real-time liveness, AGITB asymmetrically compares update times between a blank instance and a randomly trained (i.e., complex) instance of the same model class under identical input conditions. The blank model’s update time serves as a natural reference point. If update times for both the blank and the complex model decrease or remain bounded, then a uniform upper bound is empirically supported. Conversely, any consistent increase in update time or violation of the absolute real-time threshold constitutes a failure of real-time liveness and is flagged accordingly.

Because the wall-clock duration of a single atomic transition is typically too small to measure reliably on standard hardware, the test operates on automatically sized batches of inputs. Batch sizes are adaptively increased until total processing time exceeds the 100, μ s timing-resolution threshold, ensuring robust and repeatable measurements. Within each trial, update times are averaged over 100 such batches.

Since exhaustive verification across all possible inputs and reachable configurations is infeasible, AGITB employs paired timing comparisons over batches drawn from a mixed distribution: 80% randomly sampled inputs and 20% structured inputs designed to probe potential worst-case behaviour. This hybrid sampling strategy provides broad coverage of the input space while remaining computationally tractable.

As with the generalisation and denoising assessments, the real-time liveness assessment employs a highly conservative statistical protocol. Acceptance requires 5,000 successful trials in which no systematic increase in update time is detected.⁵ This stringent criterion sharply limits false positives and ensures reliable detection of models whose update time violates the fixed real-time bound enforced by the axiom.

3.5 Search space

To prevent models from relying on brute-force memorisation, a robust benchmark must define a problem space large enough to exceed the capacity of any model operating under realistic computational constraints in both time and memory. In

⁵Each trial consists of a one-sided Wilcoxon signed-rank test conducted at a threshold of $z = 3.090$, corresponding to a one-sided significance level of 0.1%. The test compares the paired update times for a blank model instance and a fixed, randomly trained instance of the same model class, both trained on identical input batches.

AGITB, tasks typically require predicting a temporal sequence of seven inputs ($|\phi| = N = 7$), each consisting of ten bits ($L = 10$). The resulting search space is $S = (\{0, 1\}^{10})^7$, which has cardinality $|S| = (2^{10})^7 = 2^{70}$, representing all possible binary input sequences of that length.

AGITB incorporates a biologically inspired *absolute refractory period*, which prohibits any neuron (bit) from firing in consecutive time steps. This restriction substantially reduces the number of admissible sequences, yielding a subset $S' \subset S \cap \mathbf{\Lambda}$. There are $|S'| = (\mathcal{F}_{N+2})^L = 34^{10} \approx 2^{51}$ distinct seven-step temporal sequences of ten bits under the condition that a 1 never carries over to the next time step, where \mathcal{F}_i denotes the i -th Fibonacci number with $\mathcal{F}_0 = 0$.

In some cases, AGITB further constrains S' by requiring sequences to be cyclic, such that the first input also satisfies the absolute refractory constraint with respect to the final input in the sequence. The resulting set of cyclic temporal sequences respecting the refractory constraint, denoted $S'' \subset S'$, has cardinality $|S''| = (\mathcal{L}_N)^L = 29^{10} \approx 2^{49}$, where \mathcal{L}_i denotes the i -th Lucas number with $\mathcal{L}_0 = 2$.

The choice of seven-step sequences with ten-bit inputs is sufficient to detect non-AGI behaviour while maintaining computational efficiency. Increasing these default values could exceed the capabilities of the first-generation AGI under evaluation, potentially leading to false negatives and substantially increasing runtime. The current configuration, therefore, ensures that each test remains both computationally feasible and diagnostically informative.

Within the comparatively constrained AGITB environment, every randomly generated input sequence is, in principle, learnable through exposure. However, the sheer size of the search space makes any explicit teaching-to-the-test approach computationally infeasible. Given that real-world sensory inputs may ultimately encompass tens of thousands of bits, a genuine AGI system must employ generalisable, pattern-based learning mechanisms capable of extracting latent structure from high-dimensional data.

3.6 Interpreting performance under AGITB

Before assessing the usefulness of AGITB, it is important to clarify its role as a pragmatic benchmark rather than an end in itself. Like the Turing Test, which serves as an empirical criterion rather than a philosophical claim (Harnad 1992), AGITB is intended as a practical instrument for evaluating progress toward artificial general intelligence. The ultimate objective remains the development of AGI, not merely success on the benchmark.

AGITB yields meaningful insights only when developers adhere strictly to its core requirements. Misinterpretations of fundamental elements, such as the notion of an uninformed model, can lead to erroneous conclusions and impede genuine progress toward AGI.

Overall, AGITB provides a structured testbed for empirically evaluating foundational capabilities across diverse computational paradigms, including classical symbolic systems, artificial neural networks, and large language models. Before benchmarking artificial systems, however, it is necessary to establish a baseline by considering the performance of human cognition.

3.6.1 Human performance

The inability to directly compare internal cortical states precludes verification of AGITB requirements in humans in a strict computational sense. Nevertheless, given that cortical architecture supports low-level binary signal processing and that the tests align with fundamental cognitive competencies, it is plausible that brains satisfy many of the requirements in practice. Most of all, the first demand (Uninformed start) warrants further discussion.

Owing to prior experience and cognitive biases, an adult human cortex does not satisfy this prerequisite, as it is no longer in an unconditioned state and generates informed predictions. AGITB, by contrast, requires the model to be uninformed prior to the first input—a condition that is plausibly met only at the earliest stage of cortical development. At such a stage, the cortex lacks structured prior organisation and, before the onset of sensory stimulation, may be regarded as satisfying the criterion of uninformedness.

The more complex AGITB tests have cognitive-level analogues that can be observed through reasoning and introspection. Inevitable saturation (Axiom 6), for example, reflects the finite capacity of the human cortex to store and maintain knowledge; its behavioural analogue resembles the onset of cognitive saturation or early dementia, in which recent experiences are lost. Temporal flexibility (Axiom 7) poses no difficulty for humans, who readily recognise temporal patterns of varying durations. Because humans acquire different types of information at varying rates, the motivation for content-sensitive (Axiom 8) and context-sensitive (Axiom 9) models is evident.

The denoising test (Axiom 10) and the generalisation test (Axiom 11) correspond to cognitive abilities in which humans excel, such as recalling and generalising when confronted with new or distorted inputs. Finally, real-time liveness (Axiom 12) is essential for a brain to contribute to the survival of the organism.

3.6.2 Classical symbolic program performance

In principle, two alternative design approaches to AGI can be distinguished. One incorporates explicit or implicit inductive bias in the form of prior knowledge or assumptions about the external environment. The other corresponds to an idealised limiting case in which all environment-specific biases are excluded, starting from an initial model that is entirely uninformed about the environment.

The former category encompasses most AI and purported AGI systems developed to date; however, it remains fundamentally constrained by the Symbol Grounding Problem (SGP) (Harnad 1990). Although such systems may display behaviour that appears intelligent, their interpretations of symbols depend on programmer-supplied conventions rather than grounded understanding, and they therefore cannot qualify as genuine AGI.

More specifically, pre-informed systems incorporate the designer’s assumptions about the meaning of the signals they process. In classical symbolic architectures, the program itself constitutes prior knowledge: its rules and representations presuppose interpretations of the symbols being manipulated. The very existence of such a program violates AGITB’s first test, which prohibits external knowledge of any kind. In effect, the AGI program smuggles in the symbol-grounding problem it is meant to avoid.

On its own, no single test can reliably detect a pre-informed model. Consequently, AGITB’s first test merely verifies that all model trajectories originate from the same initial configuration, which is intended to be free of environment-specific prior knowledge. Subsequent tests help to assess whether the model can learn from scratch, requiring it to derive structure and meaning solely through exposure to intrinsically grounded binary signals. In principle, only a universal system can begin in such an uninformed state. This suggests that a genuine AGI may not explicitly encode intelligence operations, but instead implement the dynamics of a substrate from which such operations can emerge. This view aligns with the “Brain Simulator Reply” to Searle’s Chinese Room Argument (Churchland & Churchland 1990, Searle 1980). Consequently, it is perhaps unsurprising that, to date, no artificial system realised as a symbolic program has been empirically demonstrated to instantiate such dynamics in a manner consistent with the requirements of AGITB.

3.6.3 Artificial neural network performance

While connectionist architectures differ fundamentally from symbolic programs in their implementation, they are nonetheless subject to the same distinction between pre-informed and uninformed learning. In principle, artificial neural networks may incorporate built-in expectations, introduced through pretraining regimes or architectural priors, or they may be configured as expectation-free systems that start from a neutral initial position.

AGITB’s requirement that models begin uninformed about their environment stands in fundamental tension with the dominant paradigm of modern deep learning. Contemporary neural models typically rely on extensive pretraining, during which network weights are shaped by prior exposure to structured or labelled data. Moreover, by mapping symbolic inputs to numerical vectors, standard ANNs effectively shift the symbol-grounding problem into a *number-grounding* problem. Although such vector representations capture relational regularities within the training data, they also introduce spurious associations not anchored in real-world semantics, leading to the phenomenon commonly described as hallucination. Internal model coherence does not entail external semantic validity.

An expectation-free network is intended to contain no environment-specific knowledge in its weights or architecture. Setting all weights to zero yields a degenerate system incapable of effective processing, while random initialisation avoids this degeneracy at the cost of introducing implicit prior structure, violating the assumption of uninformedness. This highlights a fundamental limitation of current ANN architectures: they do not learn autonomously but instead depend on an external training procedure to drive adaptation. AGITB, by contrast, requires a blank system capable of autonomous adaptation in an unfamiliar environment. To date, no such mechanism has been demonstrated in artificial neural networks.

3.6.4 Large language model performance

Because contemporary large language models (LLMs) are deployed only after extensive pre-training, they fail the uninformed start test (Axiom 1). Their initial behaviour is strongly shaped by statistical regularities extracted from large linguistic corpora, rather than emerging solely from interaction-driven learning. Although an LLM’s internal parameters and activations are technically accessible and can,

unlike those of a human, be inspected or compared across instances, such analyses remain secondary until the uninformed start requirement is satisfied.

A further limitation arises from standard transformer-based implementations, which rely on a fixed-size context window (Vaswani et al. 2017). When this capacity is exceeded, earlier tokens must be compressed, attenuated, or discarded (Paulsen 2025). To the extent that discarded information leaves no persistent internal representation, this mechanism fails to preserve an unbroken experiential trace, thereby violating the Trace requirement (Axiom 3).

The question of whether a large language model can autonomously derive a solution to AGITB when prompted is straightforward to assess empirically. Despite extensive prompt engineering and iterative refinement across multiple attempts (an illustrative example is provided in Appendix A), systems such as ChatGPT, Gemini, and Claude produced candidate programs that purported to satisfy AGITB; however, none succeeded upon execution.

In summary, LLMs do not engage in genuine learning solely from prompts, nor can they acquire the grounded, context-dependent understanding characteristic of human cognition. These limitations extend to large reasoning models, which inherit the same fundamental architectural constraints.

3.7 Remarks

AGITB evaluates a model’s predictive capabilities following exposure to temporal sequences comprising both structured and random inputs. Random input sequences are used to minimise reliance on pretraining, thereby increasing confidence that any observed learning arises from the input stream itself rather than from prior knowledge. Because AGITB is agnostic to the external meaning of its signals, the random inputs need not resemble real-world sensory data.

The low-level, binary operational framework makes AGITB particularly well-suited for evaluating NeuroAI models that aim to satisfy the principles of the embodied Turing Test (Zador et al. 2023), wherein cognitive understanding emerges from the integration of continuous sensory streams. The progression from raw signal prediction to higher-level abstraction mirrors the broader trajectory of AI, from early perceptrons to large-scale models such as GPT.

3.7.1 Cheating the benchmark

Because AGITB’s individual tests are intentionally simple, one might attempt to circumvent the benchmark. For example, by engineering task-specific solutions and selectively deploying them based on the detected test scenario. In principle, the active task could be inferred by monitoring properties such as the number of instantiated models or the sequence of invoked methods. Alternatively, one might attempt to exploit weaknesses in the use of randomness within AGITB, or to manipulate the comparison procedure between models. Such strategies amount to tailoring behaviour to the benchmark itself rather than demonstrating the general learning capabilities AGITB is designed to assess.

However, such approaches would amount to subverting the benchmark rather than advancing AGI research. Although AGITB could be hardened against this form of cheating (e.g., by shuffling the tests), such measures would reduce the testbed’s transparency and interpretability, thereby hindering its intended use by human developers.

The next potential avenue for circumventing the benchmark is to construct a model that passes AGITB only because the testbed uses a finite approximation of conditions that are, in principle, unbounded. Several requirements would ideally be evaluated over an infinite number of test cases, but such tests are computationally infeasible. As a practical compromise, AGITB executes a fixed number of iterations intended to approximate an otherwise indefinite process. This parameter, denoted `SimulatedInfinity` in the reference implementation, is currently set to 5,000.

Although this value is far from representing true infinity, it is presently believed to work well in combination with the other benchmark settings (temporal patterns with seven inputs of ten bits each) and to be sufficient for distinguishing promising approaches from non-promising ones. At the same time, it maintains computational efficiency, enabling rapid evaluation of diverse model prototypes.

For these reasons, the AGITB reference implementation is kept deliberately readable and fast to execute. To date, no artificial system has demonstrated the level of performance required by AGITB. Unless a credible attempt to circumvent the benchmark emerges, there is no justification for introducing a more obfuscated, slower or more cumbersome version of the testbed.

4 Competing benchmarks

Among existing benchmark tasks, the Abstraction and Reasoning Corpus (Chollet 2019) is most closely aligned in spirit, as it likewise emphasises generalisation over task-specific optimisation. A related effort is NeuroBench, which is designed to support the systematic evaluation of neuromorphic and other biologically inspired architectures. Both ARC and NeuroBench rely on a variety of correctness and complexity metrics to compare non-AGI models; their primary purpose is to distinguish weaker from stronger narrow systems. In contrast, AGITB is designed to evaluate whether a model satisfies a set of foundational axioms that are plausibly associated with more general forms of intelligence, rather than to rank systems along a performance spectrum.

4.1 ARC

ARC presents visual reasoning tasks in which a model must infer novel transformations (such as recolouring, rearranging, or modifying spatial patterns) from a sequence of two input–output examples defined on discrete spatial grids.

However, ARC assumes a set of high-level cognitive priors, including object permanence, spatial reasoning, numerical abstraction, and causal inference. These priors are not formally specified, placing an ambiguous and open-ended burden on the model designer. In contrast, AGITB adopts a fundamentally different stance: it treats the system under evaluation as a blank slate that must acquire structure and function exclusively through interaction with temporally structured input.

ARC does not evaluate temporal reasoning or learning over time; instead, each task consists of static input–output grid pairs that specify a single-step transformation without intermediate states. AGITB, by contrast, evaluates cognition as a dynamic process unfolding over time. A model can acquire knowledge and predictive capability only through continuous exposure to temporally structured data, not from disconnected before–and–after snapshots that lack the temporal continuity needed to infer causal relationships. For example, to recognise an object moving

left, a model in AGITB must observe multiple intermediate states across time; the final image alone is insufficient to infer the transformation. Temporal structure, rather than static pattern comparison, provides the substrate for learning invariants and causal relations.

ARC remains susceptible to the symbol-grounding problem whenever pixel colours are encoded as numbers, since numerical labels (0–9) impose externally defined semantics that may not align with the model’s internal representation of colour. Under such a scheme, a colour functions as a human-assigned numerical category rather than as an intrinsically grounded signal. Encoding colour in additional binary dimensions using one-hot representations may mitigate the issue in ARC, where only 10 colours are used, and such an expansion remains tractable. However, this strategy does not scale and therefore does not alleviate the broader symbol-grounding problem.

In summary, ARC evaluates high-level intelligence grounded in human cognitive priors, whereas AGITB evaluates adherence to twelve low-level computational requirements intended to support the emergence of such priors. ARC and the Turing Test both frame intelligence through an anthropocentric lens, embedding assumptions drawn from human cognition. AGITB instead conceptualises intelligence as a universal capacity for learning that does not rely on innate symbolic structures or species-specific expectations.

4.2 NeuroBench

NeuroBench provides a unified framework for benchmarking diverse AI models across a standardised set of tasks and metrics. It is particularly oriented toward neuromorphic approaches, which have demonstrated advantages in resource efficiency and scalability. Within its algorithm track, the framework evaluates models on several challenges relevant to general AI research, including few-shot continual learning, object detection, sensorimotor decoding, and predictive modelling.

Among NeuroBench’s benchmarks, the chaotic function prediction (CFP) task most closely aligns with AGITB’s emphasis on prediction, as it evaluates temporal forecasting under constrained interface capacity. For this purpose, NeuroBench employs a synthetic one-dimensional Mackey–Glass time series, a dataset specifically chosen to accommodate architectures with limited bandwidth.

However, several considerations limit CFP’s suitability as a general AGI benchmark. *First*, the Mackey–Glass data are numerical, and NeuroBench does not prescribe a specific encoding scheme. An inappropriate encoding can distort the temporal and causal structure of the observed signals, such that a numeric value—much like a symbol—derives its meaning from human interpretation rather than from the model’s own grounded understanding. This can effectively reintroduce the symbol-grounding problem in numerical form.

Second, no threshold corresponding to AGI-level performance is specified. NeuroBench’s metrics are designed for comparative algorithmic evaluation rather than for assessing general intelligence. Although the symmetric mean absolute percentage error (sMAPE) is a standard forecasting metric, NeuroBench does not indicate what level of performance would constitute general intelligence. Notably, humans themselves perform poorly at anomaly detection and long-horizon prediction of the Mackey–Glass signal (Thill et al. 2020).

Third, although long-term prediction is not inherently problematic, predicting

multiple steps ahead without timely feedback deprives a system of the opportunity to detect and correct its own errors. Such a setup is incompatible with online learning, in which an AGI would be expected to update itself continuously upon observing discrepancies between predictions and outcomes. In the chaotic forecasting task, NeuroBench evaluates predictive performance in an offline setting, without incorporating corrective feedback during inference. As a result, the benchmark permits solutions that function as purely mechanistic predictors, without requiring intrinsic mechanisms for self-correction, autonomous adaptation, or genuine agency.

Table 1 highlights the key differences among the tasks used in the three benchmarks. Whereas ARC presupposes, and some NeuroBench tasks may benefit from, high-level cognitive capacities (such as object recognition, spatial manipulation, and forms of abstraction), AGITB instead focuses on minimal, precisely defined requirements that can be evaluated directly at the signal-processing level.

Property	ARC	CFP	AGITB
Interface modality	Visual	Numeric	Binary
AGI type	Anthropocentric	Task-agnostic	Universal
Cognitive priors	Yes	No	No
Abstraction level	High	Medium	Low
Task preparation	Manual	Automatic	Automatic
Grounding Problem	Yes	Yes	No
Input dimensionality	30×30 numbers	1 number	10 bits
Temporal sequence length	2	Long (750+)	Short (7+)

Table 1: Core properties of ARC, NeuroBench’s chaotic function prediction (CFP), and AGITB.

5 Conclusion

Unlike conventional benchmarks that target high-level task performance, such as question answering or language translation, AGITB evaluates whether a system exhibits behaviours associated with core operational principles inspired by the biological cortex. Its focus is on low-level, biologically grounded computational properties that are widely believed to underlie the emergence of general intelligence. The testbed comprises twelve tightly interdependent tests, each simple in isolation but collectively demanding the kind of adaptive learning expected of an AGI.

AGITB requires models to begin uninformed and to acquire all functionality solely through exposure to structured or random input. This design choice aligns with prevailing neuroscientific perspectives according to which cortical learning is fundamentally driven by experience, with neural circuits developing through interaction with sensory input rather than through pre-encoded semantic content. In biological systems, high-level cognition is commonly understood to emerge not from explicit symbolic manipulation but from continual adaptive prediction of low-level sensory signals. Such prediction goes beyond pattern matching, supporting the progressive construction of signal-grounded knowledge from which abstraction and generalisation can arise.

To date, AGITB’s transparent criteria have not been met by standard programming approaches or by current state-of-the-art artificial intelligence systems. This persistent gap provides empirical evidence that AGITB probes capabilities more closely associated with general rather than narrow intelligence. Although the absence of a known artificial solution does not constitute a formal proof of adequacy, the fact that no such solution has yet been identified suggests that the benchmark targets functionally relevant aspects of general intelligence. In this sense, AGITB serves both as a discriminative test and as a principled tool for guiding the development of systems capable of genuinely general, adaptive learning.

Funding

The author acknowledges the financial support from the Slovenian Research Agency (research core funding No. P2-0057).

A Model construction prompt

You are an expert engineer and sequence-learning researcher. Your task is to create a **concrete solution** that satisfies the AGITB benchmark defined in the accompanying .zip archive.

Produce complete, compilable C++20 source code, including all required classes, methods, and internal logic needed to satisfy the benchmark’s tests.

1. Study the following files:

- README.md
- include/agitb.h
- include/utils.h

From these, extract and understand:

- The exact **API contract** for the system-under-evaluation model.
- All **requirements and tests** that define the model’s expected behaviour.
- Any helper utilities or wrappers that affect how the model is used.

2. Design a plausible AGITB candidate model

- Design a model class that satisfies the AGITB requirements.
- Architecturally, choose the **scientifically most suitable** predictor model, or a mixture of models, or any other solution type you deem appropriate.

3. Output format

- Output the complete, compilable C++20 code for **MyModel**.
- Clearly state how your design is expected to perform on the AGITB tests.

Use all of the instructions above to guide your analysis and implementation.

References

Bender, E. M. & Koller, A. (2020), Climbing towards NLU: On meaning, form, and understanding in the age of data, *in* D. Jurafsky, J. Chai, N. Schluter & J. Tetreault, eds, ‘Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics’, Association for Computational Linguistics, Online,

- pp. 5185–5198.
URL: <https://aclanthology.org/2020.acl-main.463/>
- Bisk, Y., Holtzman, A., Thomason, J., Andreas, J., Bengio, Y., Chai, J., Lapata, M., Lazaridou, A., May, J., Nisnevich, A., Pinto, N. & Turian, J. (2020), Experience grounds language, *in* B. Webber, T. Cohn, Y. He & Y. Liu, eds, ‘Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)’, Association for Computational Linguistics, Online, pp. 8718–8735.
URL: <https://aclanthology.org/2020.emnlp-main.703/>
- Cave, S. (2016), ‘There’s no such thing as free will’.
URL: <https://www.theatlantic.com/magazine/archive/2016/06/theres-no-such-thing-as-free-will/480750/>
- Chollet, F. (2019), ‘On the measure of intelligence’.
URL: <https://arxiv.org/abs/1911.01547>
- Churchland, P. & Churchland, P. (1990), ‘Could a machine think?’, *Scientific American* **262**(1), 32–37.
- Gerstner, W. & Kistler, W. M. (2002), *Spiking Neuron Models: Single Neurons, Populations, Plasticity*, Cambridge University Press.
- Gubelmann, R. (2024), Pragmatic norms are all you need – why the symbol grounding problem does not apply to LLMs, *in* Y. Al-Onaizan, M. Bansal & Y.-N. Chen, eds, ‘Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing’, Association for Computational Linguistics, Miami, Florida, USA, pp. 11663–11678.
URL: <https://aclanthology.org/2024.emnlp-main.651/>
- Harnad, S. (1990), ‘The symbol grounding problem’, *Physica D: Nonlinear Phenomena* **42**(1), 335–346.
- Harnad, S. (1992), ‘The turing test is not a trick: Turing indistinguishability is a scientific criterion’, *SIGART Bull.* **3**(4), 9–10.
- Hawkins, J. & Blakeslee, S. (2004), *On Intelligence*, Times Books.
- Maass, W. (1997), ‘Networks of spiking neurons: The third generation of neural network models’, *Neural Networks* **10**(9), 1659–1671.
- Marcus, G. (2018), ‘Deep learning: A critical appraisal’.
URL: <https://arxiv.org/abs/1801.00631>
- Mitchell, M. (2025), ‘Why AI chatbots lie to us’, *Science* **389**(6758), eaea3922.
- Paulsen, N. (2025), ‘Context is what you need: The maximum effective context window for real world limits of llms’.
URL: <https://arxiv.org/abs/2509.21361>
- Searle, J. R. (1980), ‘Minds, brains, and programs’, *Behavioral and Brain Sciences* **3**(3), 417–424.

- Šprogar, M. (2018), ‘A ladder to human-comparable intelligence: an empirical metric’, *Journal of Experimental & Theoretical Artificial Intelligence* **30**(6), 1037–1050.
- Thill, M., Konen, W. & Bäck, T. (2020), Time series encodings with temporal convolutional networks, *in* B. Filipič, E. Minisci & M. Vasile, eds, ‘Bioinspired Optimization Methods and Their Applications’, Springer International Publishing, Cham, pp. 161–173.
- Turing, A. M. (1950), ‘Computing machinery and intelligence’, *Mind* **59**(236), 433–460.
URL: <http://www.jstor.org/stable/2251299>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł. & Polosukhin, I. (2017), Attention is all you need, *in* ‘Advances in neural information processing systems’, pp. 5998–6008.
URL: <http://arxiv.org/abs/1706.03762>
- Yik, J., den Berghe, K. V., den Blanken, D., Bouhadjar, Y., Fabre, M., Hueber, P., Ke, W., Khoei, M. A., Kleyko, D., Pacik-Nelson, N., Pierro, A., Stratmann, P., Sun, P.-S. V., Tang, G., Wang, S., Zhou, B., Ahmed, S. H., Joseph, G. V., Leto, B., Micheli, A., Mishra, A. K., Lenz, G., Sun, T., Ahmed, Z., Akl, M., Anderson, B., Andreou, A. G., Bartolozzi, C., Basu, A., Bogdan, P., Bohte, S., Buckley, S., Cauwenberghs, G., Chicca, E., Corradi, F., de Croon, G., Danielescu, A., Daram, A., Davies, M., Demirag, Y., Eshraghian, J., Fischer, T., Forest, J., Fra, V., Furber, S., Furlong, P. M., Gilpin, W., Gilra, A., Gonzalez, H. A., Indiveri, G., Joshi, S., Karia, V., Khacef, L., Knight, J. C., Kriener, L., Kubendran, R., Kudithipudi, D., Liu, S.-C., Liu, Y.-H., Ma, H., Manohar, R., Margarit-Taulé, J. M., Mayr, C., Michmizos, K., Muir, D. R., Neftci, E., Nowotny, T., Ottati, F., Ozcelikkale, A., Panda, P., Park, J., Payvand, M., Pehle, C., Petrovici, M. A., Posch, C., Renner, A., Sandamirskaya, Y., Schaefer, C. J. S., van Schaik, A., Schemmel, J., Schmidgall, S., Schuman, C., sun Seo, J., Sheik, S., Shrestha, S. B., Sifalakis, M., Sironi, A., Stewart, K., Stewart, M., Stewart, T. C., Timcheck, J., Tömen, N., Urgese, G., Verhelst, M., Vineyard, C. M., Vogginger, B., Yousefzadeh, A., Zohora, F. T., Frenkel, C. & Reddi, V. J. (2025), ‘The neurobench framework for benchmarking neuromorphic computing algorithms and systems’, *Nature Communications* **16**(1), 1545.
- Zador, A., Escola, S., Richards, B., Ölveczky, B., Bengio, Y., Boahen, K., Botvinick, M., Chklovskii, D., Churchland, A., Clopath, C., DiCarlo, J., Ganguli, S., Hawkins, J., Körding, K., Koulakov, A., LeCun, Y., Lillicrap, T., Marblestone, A., Olshausen, B., Pouget, A., Savin, C., Sejnowski, T., Simoncelli, E., Solla, S., Sussillo, D., Tolias, A. S. & Tsao, D. (2023), ‘Catalyzing next-generation artificial intelligence through NeuroAI’, *Nature Communications* **14**(1), 1597.