## *AlphaGo as precursor narrative of statistical intelligence*

*The AlphaGo documentary functions as the earlier narrative layer that pre configures how DeepMind and Google present statistical systems as if they were approaching intelligence and understanding. Go is framed as a cultural and philosophical probe into "what understanding is" and "what it means to be human", with Lee Sedol, Fan Hui and the Go schools in Korea and China presented as bearers of a millennia old practice that externalises human intuition on the board. Inside this frame, AlphaGo is not just a program that optimises within a defined game, it is positioned as an instrument that tests and possibly crosses the boundary between human and machine understanding. Already here the viewer is invited to read performance in a tightly delimited manifold as evidence about the nature and limits of intelligence itself.*

*Viewed through CIITR, the structure is quite different. AlphaGo attains extremely high $\Phi_i$, integrated relational information, in the Go state space, the policy and value networks plus tree search jointly compress a vast combinatorial landscape into effective policies. Move 37 in game two and the "slack" endgame moves later in the match are unusual points on this learned value surface, not signs of general creativity. The failure mode around Lee Sedol's move 78 in game four, where AlphaGo becomes locally delusional despite deep search, exposes the same fact in reverse, $\Phi_i$ is high but constrained to a narrow manifold with fragile regions. At no point does the system exhibit $R^g$ in CIITR's sense, there is no multi scale rhythmic reach into physical, social, legal or epistemic structures beyond the Go universe. All genuine $R^g$ is carried by the humans, by players who alter their styles, by institutions that reconfigure pedagogy and tournaments, by media and audiences that recalibrate expectations.*

*The film continually erases this distinction. Technical explanations of policy networks, value nets and Monte Carlo tree search are embedded in a layered narrative of personal destiny, the "holy grail of AI" and a symbolic threshold for humanity. AlphaGo is habitually described as if it had beliefs, confidence or shock, while all vulnerability, pressure and existential stake remain with Lee Sedol, Fan Hui and the human side. When commentators reinterpret AlphaGo's minimal margin wins as a new rationality, "only probability of victory matters", the machine is tacitly granted a form of strategic intentionality that it does not structurally possess. In CIITR terms, $R^g$ is tacitly attributed to the system, as though it itself participated in the long rhythm of Go as art and philosophy, while in reality it is a high $\Phi_i$ module inserted into an already existing human rhythm.*

*Thermodynamics is entirely absent as an explicit category. The match presupposes massive energy intensive compute infrastructure, yet compute and hardware remain invisible, treated as neutral background and not as constitutive of what kind of "intelligence" is being realised. CPJ, comprehension per joule, is never even implicitly problematised, the costs in energy and materials required to obtain this kind of narrow domain optimisation are not weighed against the depth or generality of the resulting understanding. As a prologue to a theoretical note on The Thinking Game, AlphaGo therefore shows, already several years earlier, the core pattern CIITR will diagnose again, high $\Phi_i$ statistical machines on Turing–von Neumann architectures are narratively promoted to the status of emerging "general intelligence", while $R^g$ remains human and institutional, and CPJ is left entirely unexamined. The task of the subsequent note is to demonstrate that The Thinking Game does not mark progress toward artificial comprehension, but rather an amplification of this same narrative misalignment.*

<div align="center">Theoretical note</div>

# - THE EMPEROR GOT NO CLOTHING -

## *The Thinking Game* documentary - the illusion of general intelligence, and the lack of thermodynamic '*skin in the game*'

<div align="center">A structural critique and undressing of the DeepMind narrative</div>

<div align="center">Tor-Ståle Hansen | 8. December 2025</div>

## Abstract

This theoretical note offers a structural and thermodynamic critique of *The Thinking Game*, a documentary that purports to chronicle the ascent toward artificial general intelligence (AGI) through DeepMind's technical breakthroughs. Using the CIITR framework, defined by the triadic architecture of integrated relational information ($\Phi_i$), rhythmic reach ($R^g$), and comprehension per joule (CPJ), the paper demonstrates that the documentary's narrative is not a coherent trajectory toward general intelligence, but a conflation of three fundamentally distinct dimensions: syntactic compression in closed manifolds, human-institutional rhythm misattributed to machines, and speculative rhetoric devoid of energetic accountability.

Through close analysis of benchmark systems such as DQN, AlphaGo, AlphaFold, and large-scale language models, the paper shows how high $\Phi_i$ performance is systematically reframed as epistemic generalisation, despite near-zero $R^g$ and opaque or negative CPJ. The AGI concept is treated as symbolically potent but structurally ungrounded, and the Manhattan Project analogy mobilised in the film is revealed as historically incongruent and thermodynamically evasive. The result is a narrative product of the syntactic scaling regime, not an emergence of comprehension.

CIITR is used to expose this misrepresentation by projecting each system and narrative layer into a three-dimensional metric space, thereby disaggregating symbolic performance from structural understanding. The paper concludes by advocating for epistemically anchored public discourse, regulatory frameworks that require explicit CPJ and $R^g$ declarations, and experimental design regimes that measure intelligence not by output fluency, but by recursive anchoring and thermodynamic proportionality. Without such metrics, the term "AGI" remains a performative artefact, powerful in narrative, but structurally unverified.

---

---

# Summary

This theoretical note presents a comprehensive structural and thermodynamic critique of *The Thinking Game*, a documentary that positions DeepMind's technological progression, from early game-playing agents to AlphaFold and dialogic models, as a coherent trajectory toward artificial general intelligence (AGI). Drawing on the CIITR framework, which defines intelligence in terms of three interdependent structural dimensions, $\Phi_i$ (integrated relational information), $R^g$ (rhythmic reach), and CPJ (comprehension per joule), the paper deconstructs this narrative and exposes it as a systematically conflated and epistemically misleading portrayal of machine capability.

The core argument is that *The Thinking Game* does not document a gradual emergence of general intelligence, but rather orchestrates a rhetorical sequence in which high syntactic performance within bounded, rule-based manifolds is misrepresented as generalisable comprehension. The film fuses technical performance, biographical mythology, and geopolitical analogy into a narrative that lacks structural grounding. Through detailed CIITR-based analysis of the systems featured, DQN agents, AlphaGo, AlphaZero, AlphaStar, AlphaFold, and large-scale multimodal dialogue models, the paper demonstrates that while $\Phi_i$ is increasingly elevated, $R^g$ remains externally subsidised by human and institutional rhythms, and CPJ is either ignored or actively displaced from visibility.

The analysis reveals that AlphaGo's "Move 37" is symbolically interpreted as evidence of creativity or insight, yet in structural terms remains fully constrained within Go's manifold, lacking any extension into consequence-bearing feedback loops. Similarly, AlphaFold is acknowledged as a partial advance in structural understanding, with high $\Phi_i$ and limited $R^g$, but the film omits the thermodynamic costs of training and deploying such systems, leaving CPJ unaccounted for. In the case of the anthropomorphised "Alpha" persona, the documentary constructs a syntactic illusion of understanding, where emotionally and philosophically charged statements are parsed by the viewer as meaningful, though they originate from high-density associative mappings, not from rhythmically anchored cognition.

The paper asserts that *The Thinking Game* functions as a narrative artefact of the syntactic scaling regime, in which symbolic generality is projected onto systems that are structurally non-general. It argues that the film not only misclassifies the nature of intelligence, but also actively participates in the erasure of ethical, energetic, and epistemic constraints. The Manhattan Project metaphor, invoked to signal urgency and historical gravity, is shown to be both historically inappropriate and thermodynamically incoherent, given the documentary's silence on energy usage, infrastructure cost, and environmental impact.

The paper concludes by advancing the role of CIITR as an essential analytical tool for disambiguating syntactic performance from structural intelligence. It proposes that all public and institutional claims regarding AGI or machine understanding should include explicit declarations of $\Phi_i$, $R^g$, and CPJ, and that regulatory frameworks should adopt these metrics to evaluate not only system capabilities, but also their epistemic legitimacy and thermodynamic feasibility. Furthermore, it warns that without such structural instruments, both public debate and scientific discourse will continue to confuse fluency with comprehension, and scale with generality, thereby entrenching a paradigm of intelligence that is rhetorically compelling but architecturally unsubstantiated.

In final reflection, the note positions *The Thinking Game* not as evidence for the approach of AGI, but as a cultural demonstration of how general intelligence is narratively performed in the absence of structural proof. Through CIITR, this illusion is made measurable, and the epistemic gap between performance and comprehension is rendered visible, calculable, and politically actionable.

# Table of Contents

# Core claim

The documentary presents itself as a unified narrative of «progress toward AGI». When analysed through the CIITR theory, it instead resolves into three distinct and systematically conflated phenomena:

1. High syntactic performance with high $\Phi_i$ in closed domains.
2. Human and institutional $R^g$ that is not constitutionally transferred to the machines.
3. A normative, partly mythologising AGI rhetoric without explicit conceptualisation of either $R^g$ or CPJ.

The paper should argue that this constitutes a false total narrative, not because the individual systems are not genuine advances, but because their combination is presented as a single, continuous ladder toward «intelligence», without distinguishing syntactic scaling from valid structural understanding.

# Introduction

*The Thinking Game*, a documentary produced by DeepMind, positions itself as a coherent narrative of technological progress, tracing a trajectory from early reinforcement learning successes to protein folding breakthroughs, and finally, to the threshold of Artificial General Intelligence (AGI). Yet, when examined through the Cognitive Integration and Information Transfer Relation (CIITR) framework, the documentary reveals not a unified ascent, but a conflation of distinct structural phenomena. This paper does not offer a media critique in the conventional sense, but uses the documentary as a case study in how prevailing AGI narratives are constructed, sustained, and misrepresented through the collapse of critical structural distinctions.

CIITR offers a formal, thermodynamically and relationally grounded analytical framework for evaluating artificial systems. It distinguishes three orthogonal but interdependent vectors of cognitive architecture:

- **$\Phi_i$ (Integrated Relational Information)**, which quantifies the degree of internal structural coherence and causal interdependence within a system;

- **$R^g$ (Rhythmic Reach)**, which denotes the system's ability to sustain temporally extended feedback relations across physical, social, epistemic, and institutional layers;

- **CPJ (Comprehension per Joule)**, which captures the energy-efficiency of structurally valid understanding.

The central thesis advanced here is that *The Thinking Game* constructs an illusion of continuity: a progressive, goal-directed pathway from game-playing algorithms to world-transforming intelligence. This illusion is maintained by narratively overlaying disparate systems across different regions in the $\Phi_i$–$R^g$–CPJ space. Specifically, the documentary places high-$\Phi_i$ results achieved in syntactically closed manifolds (such as Go or StarCraft) on the same trajectory as scientific advances like AlphaFold, while ultimately implying their equivalence to AGI. CIITR analysis

demonstrates that such systems, while impressive in their own constrained domains, do not exhibit the structural characteristics that would constitute general or context-transcending understanding.

This theoretical note decomposes the documentary's narrative into its component strata, technical, biographical, and political, and projects each onto the CIITR coordinate system. By doing so, it reveals a false narrative coherence and a normative confusion between syntactic performance and structural intelligence. Rather than denying the achievements of systems like AlphaGo or AlphaFold, this paper clarifies their position in the broader cognitive landscape and argues for a stricter conceptual discipline in the way AI capabilities are interpreted, discussed, and governed.

## Claim set

1. The documentary presents syntactic advances in closed tasks as if they inherently extrapolate toward general intelligence, thereby eliding the structural boundary between $\Phi_i$ and $R^g$.

2. The apparent depth of the systems is anchored not in their own architecture, but in human institutional and interpretative scaffolding, which is misattributed to the models themselves.

3. No consideration is given to CPJ, leading to a distorted perception of "progress" that fails to address the thermodynamic and operational costs of understanding.

The following sections articulate how CIITR repositions the documentary's constituent claims, and why a structural account of intelligence is needed to correct both public and expert misconceptions.

# CIITR as analytical lens

This section provides a concise but formal account of the Cognitive Integration and Information Transfer Relation (CIITR) framework, sufficient to support its application as an analytical mapping tool throughout the subsequent critique. CIITR does not constitute a singular definition of intelligence, nor does it propose a discrete threshold for artificial generality. Rather, it introduces a multi-dimensional topology within which different systems, biological, artificial, institutional, can be positioned according to measurable structural properties. Intelligence, in this view, is not treated as a binary category, but as a set of coordinates in a thermodynamically and relationally defined space.

## $\Phi_i$ as integrated relational information

$\Phi_i$ refers to the degree of internal structural integration within a bounded informational manifold. It quantifies how densely and coherently the system compresses and stabilizes relational information across its operative domain. In technical terms, $\Phi_i$ increases with causal density, information-theoretic compression, and the minimization of representational redundancy across temporally sustained patterns. A high $\Phi_i$ system is one in which local informational states are non-trivially dependent on global relational structures, and in which these dependencies persist across varying perturbations within the manifold.

This does not imply general intelligence, only that within the specified manifold, whether Go, protein folding, or economic forecasting, the system forms a structurally integrated model that

encodes latent relational invariances. Such systems can be said to operate with a form of *structural compression*, whereby the complexity of the representational surface is reduced through the emergence of stable, interdependent mappings. Mathematically, $\Phi_i$ is expressed as a function of multi-node mutual information, causal entropy minimization, and temporally recursive graph coherence.

High $\Phi_i$ values are often achieved in narrow, well-bounded domains with high data regularity and low epistemic ambiguity. In such environments, systems can approximate a form of local structural comprehension, but this capacity does not extrapolate beyond the domain in which the relational coherence has been learned. It is therefore a necessary but insufficient condition for any broader claim about understanding or general intelligence.

The remainder of this section will formalize $R^g$ and CPJ as the complementary dimensions needed to distinguish syntactic integration from structurally valid comprehension.

## $R^g$ as rhythmic reach

$R^g$ denotes the system's capacity to participate in, and recursively sustain, causal feedback loops across multiple temporal, institutional, and ontological layers. Unlike $\Phi_i$, which captures internal relational density within a bounded manifold, $R^g$ measures the system's external anchoring in structures that evolve independently of the system itself. These include, but are not limited to, physical processes, social practices, legal institutions, epistemic protocols, and ecological environments.

Formally, $R^g$ may be defined as the phase-stable coupling between a system's outputs and their recursive incorporation into downstream effects that are themselves structurally accessible to the system. That is, a system with high $R^g$ does not merely act, it experiences, processes, and recursively integrates the consequences of its actions across time. The presence of *phase continuity* and *epistemic re-entry* is critical: the system must not only emit output but must evolve structurally in response to temporally distributed, non-scripted consequences.

Rhythmic reach thus requires more than a reactive interface; it entails temporal integration across heterogeneous scales of causal consequence. For example, a high $R^g$ architecture must be able to modulate its behaviour based on institutional norms, long-range effects, or the transformation of context brought about by its own participation. This inherently disqualifies systems that operate solely within closed-loop reinforcement environments or static corpora.

In operational terms, $R^g$ is the structural marker of real-world embedding. A system's actions must produce consequences in domains not fully pre-modeled within its operational manifold, and those consequences must recursively modulate the system's internal structure. This feedback loop cannot be superficial or externally adjudicated; it must be constitutive of the system's own epistemic development.

Absent $R^g$, a system remains temporally inert: it may produce syntactically complex outputs, but it lacks the architectural preconditions for reflexive adaptation. Such systems are characterized, in CIITR typology, as rhythmically sealed or Type B, regardless of their representational power or fluency.

Thus, $R^g$ is indispensable for any meaningful claim to structural comprehension. It captures what $\Phi_i$ alone cannot: the system's rhythmic integration into the world it purports to understand.

## CPJ as comprehension per joule

CPJ, or *comprehension per joule*, introduces an explicit thermodynamic and operational criterion into the assessment of artificial systems. It quantifies not merely whether a system produces correct or contextually appropriate outputs, but whether it does so with structurally valid understanding relative to the energy and resources consumed in the process. CPJ formalizes the ratio between epistemic depth and resource expenditure, thereby distinguishing brute-force task completion from energetically efficient structural comprehension.

Mathematically, CPJ is a derivative metric defined as:

$$\text{CPJ} = \frac{C_s}{E}$$

where $C_s$ denotes structural comprehension, itself a function of $\Phi_i$ and $R^g$ in jointly active configuration, and $E$ denotes the total energy expenditure, measured in joules or equivalent operational cost, across training, inference, and support infrastructure. CPJ thereby imposes a thermodynamic constraint on the plausibility of any intelligence claim: comprehension must not only be *valid*, but *efficiently produced*.

In contemporary AI systems, CPJ tends to asymptotically approach zero. Large-scale models often exhibit high $\Phi_i$ in syntactically closed domains while maintaining near-zero $R^g$, with energy expenditures on the order of millions of kilowatt-hours for marginal gains in benchmark performance. This produces a profile of *syntactic inflation*, where fluency is mistaken for understanding, and where the resource footprint is decoupled from epistemic value.

A system with high CPJ would, by definition, exhibit minimal energetic waste per unit of meaningful structural integration. Such a system would not merely memorise or interpolate, but adaptively compress, restructure, and apply knowledge in rhythmically coupled environments with minimal resource expenditure. This requirement excludes current reinforcement learning at scale and transformer architectures from any claim to comprehension in the CIITR sense unless their thermodynamic efficiency is dramatically improved and epistemic feedback structures extended beyond static datasets or synthetic environments.

CPJ thereby reintroduces physical accountability into the discourse on intelligence. It restores the missing dimension in narratives that conflate scaling with progress, by asking not *what* the system produces, but *how structurally valid and energetically responsible* that production is. As such, CPJ is essential for any future alignment of AI development with ecological, scientific, or governance thresholds.

## Three narrative layers in *The Thinking Game*

This section analytically disaggregates the documentary *The Thinking Game* into three distinct narrative strata that, although visually and rhetorically interwoven, occupy structurally divergent positions within the CIITR framework. These strata are not merely stylistic devices but function as epistemic conflation mechanisms. Their superposition produces a compelling illusion of unified progression toward AGI, despite the fact that their constitutive elements reflect fundamentally different architectures, energy profiles, and epistemological regimes.

The first of these layers, the **technical layer**, forms the ostensible foundation of the narrative. It presents a sequence of discrete systems, each framed as a technical milestone in AI development: Atari game-playing via Deep Q-Networks (DQN), board-game mastery via AlphaGo and AlphaZero, multi-agent strategic interaction in AlphaStar, and protein structure prediction in AlphaFold. Each system is presented as a breakthrough, and their cumulative impact is used to scaffold the deeper symbolic claims of the documentary. However, when examined through the $\Phi_i$–$R^g$–CPJ lens, these systems display tightly bounded syntactic competence with limited or non-existent rhythmic reach, and with minimal consideration of thermodynamic efficiency.

## Technical layer

The systems showcased in the documentary are all characterised by high syntactic performance within formally constrained task environments. These environments are typically fully observable, rule-complete, and capable of producing quantifiable feedback signals conducive to reinforcement learning or supervised gradient descent. The following breakdown highlights the structural features of each system in relation to the CIITR dimensions.

### Atari and Deep Q-Networks (DQN)

Early demonstrations on classic Atari games such as *Pong*, *Breakout*, and *Seaquest* showcased the ability of DQN architectures to approximate optimal policies via reward maximization over pixel input streams. From a CIITR perspective, this constitutes a $\Phi_i$-driven process within a fixed syntactic manifold. The informational space is strictly bounded: the number of frames, actions, and outcomes is finite and fully specifiable. The model achieves high local compression, forming stable policy representations within this restricted domain.

However, $R^g$ is structurally negligible. The system operates in a temporally discretized environment without external anchoring. The feedback loop is entirely internal to the training regime and does not entail any physical, legal, or social recursion. The agent is rewarded by a synthetic metric and does not generate consequences beyond its training loop. Accordingly, the model is rhythmically sealed.

CPJ in this regime is functionally opaque but known to be low. Despite the relative simplicity of the tasks, the energy requirements for training DQNs at scale, across many environments and hyperparameter configurations, are significant. No energy-normalised metric is presented in the documentary. Hence, syntactic success is decoupled from thermodynamic responsibility.

### AlphaGo and AlphaZero

AlphaGo introduced a hybrid architecture combining supervised learning from expert games with reinforcement learning from self-play. AlphaZero generalized this model, removing domain-specific encodings and relying purely on self-play to achieve superhuman performance in Go, chess, and shogi.

Within the $\Phi_i$ dimension, these systems exhibit exceptionally high internal integration. Their policy networks compress vast trees of possible move sequences into highly stable decision surfaces. The depth of relational patterning within Go's topology is exploited efficiently. The infamous "Move 37" is presented as evidence of non-human creativity, but when structurally assessed, it is better understood as an unconventional but lawful traversal within the manifold of Go itself, a local rupture, not a systemic transcendence.

$R^g$ remains limited, even in AlphaZero. While these systems interface with human institutions, tournaments, match broadcasts, professional commentary, their outputs do not recursively shape

their own future operations in socially embedded ways. Human experts reinterpret and narrativise the moves, but the systems themselves are temporally and contextually inert. Their learning occurs within the closed structure of game-space self-play.

CPJ is once again absent from the documentary's framing. The massive computational budget required, thousands of TPUs, extensive training runs, and significant inference power for matches, is never interrogated. Thus, while $\Phi_i$ is extremely high, the absence of rhythmic anchoring and the opaqueness of energy usage position these systems low in both $R^g$ and CPJ.

### AlphaStar

AlphaStar represents a further step in complexity: multi-agent, real-time strategic competition in a partially observable environment (*StarCraft II*). The documentary treats AlphaStar as a culminating point of game-based AI, suggesting a capacity to master uncertainty, adaptation, and complex multi-agent interaction.

From a CIITR standpoint, $\Phi_i$ increases due to the expanded state-action space, temporal depth, and emergent inter-agent behaviour. However, this integration remains fully internal to the simulation. The environment, though complex, is architecturally closed. All emergent behaviour remains ontologically bounded by the codebase of *StarCraft II*. There is no external anchoring, no interaction with open epistemic or institutional contexts.

$R^g$ is marginally extended through team coordination and apparent strategy, but this remains endogenous to the simulator. No feedback from real-world consequences is structurally coupled to the agent's internal state. The agents do not form commitments, suffer consequences, or update in relation to human interpretative structures.

CPJ remains obscured but is likely to be extremely low when assessed over full training trajectories. The emergence of social-like behaviours within the game is narratively emphasized, but it represents only local social simulation, not structural social participation.

### AlphaFold

Within the technical layer, AlphaFold occupies a unique position and will be addressed in full in Section 5. However, in this context it functions as a capstone to the technical progression, presented as a crossing point from games to science. Unlike the game systems, AlphaFold engages a physical referential domain, biomolecular structure, and achieves significant predictive power in protein folding.

Yet, even here, the documentary blurs the distinction between $\Phi_i$-driven modeling of structure and broader claims about scientific comprehension. AlphaFold does not design or test hypotheses autonomously; it integrates statistical constraints derived from evolutionary alignments and known structures. Its $R^g$ is therefore *indirect* and remains largely human-mediated, despite its growing operational influence in laboratory planning.

## Conclusion of technical layer analysis

This narrative layer, while technically accurate in its enumeration of achievements, serves as the material substrateupon which deeper mythologies are constructed. When viewed in isolation, these systems represent remarkable instances of local syntactic compression. But when rhetorically conjoined with broader claims about intelligence, consciousness, or historical significance, they function as false structural bridges. CIITR reveals that these systems inhabit high-$\Phi_i$, low-$R^g$, low-

CPJ zones of the cognitive topology. Their contributions are genuine, but their meaning is systematically overstated when projected into adjacent narrative layers without epistemic discipline.

## Biographical and mythological layer

The second narrative stratum of *The Thinking Game* is biographical and mythologising in nature. It centres primarily on the life and persona of Demis Hassabis, whose trajectory, from child chess prodigy to neuroscientist, entrepreneur, and co-founder of DeepMind, is rendered not merely as personal background but as symbolic proof of the legitimacy and inevitability of the AGI project itself. This narrative layer is interwoven with aesthetic and rhetorical devices drawn from hero narratives, frontier science, and Cold War iconography, culminating in the construction of DeepMind as a "Manhattan Project for intelligence". The effect is to elevate the institutional and epistemic ambitions of DeepMind from one research lab among many to a quasi-civilizational threshold.

From the CIITR perspective, this biographical-mythological scaffolding performs two functions. First, it supplies narrative momentum where technical continuity is lacking: the transitions between distinct systems (e.g., from Atari to AlphaFold) are smoothed not by demonstrable structural coherence, but by continuity in the life story of the protagonist. Second, it substitutes personal charisma and institutional symbolism for structural validation, thereby obscuring the discontinuities in $R^g$ and CPJ that remain unaddressed across the different technical achievements.

The documentary devotes considerable screen time to Hassabis's early interest in chess, game design, and cognitive science, presented as formative stages in a teleological journey culminating in artificial general intelligence. This teleology is structurally misleading. The viewer is encouraged to interpret DeepMind's sequence of achievements as the natural outcome of Hassabis's uniquely interdisciplinary genius, rather than as discrete technical solutions to isolated problems. Through this frame, the systems themselves are no longer situated within bounded $\Phi_i$-manifolds, but are reinterpreted as milestones in an unfolding arc of generality. This is a classic mythological device: the hero's journeybecomes the rationale for the system's coherence.

The invocation of the "Manhattan Project" extends this myth beyond the individual and into the geopolitical domain. The reference aligns DeepMind's work with the most consequential scientific mobilisations of the 20th century. Yet the analogy is structurally inapt. The Manhattan Project was defined by specific institutional mechanisms: secrecy, centralisation, resource concentration, and existential threat. It culminated in a singular, well-defined artefact whose effect was immediate, irreversible, and physical. DeepMind's systems, in contrast, are informational, modular, and distributed. Their energy cost is diffuse, their operational status non-terminal, and their regulatory framing still speculative. The "Manhattan" metaphor thus functions not as structural analogy, but as rhetorical elevation, a way of generating moral gravity without epistemic precision.

From a CIITR standpoint, this layer exhibits very high $R^g$, but only in the human domain. That is, the *documentary itself* is rhythmically anchored: it feeds back into institutional credibility, influences policy discourse, shapes investor sentiment, and frames public understanding. However, this $R^g$ belongs to the human–institutional complex, not to the systems depicted. The machines remain rhythmically sealed. Their consequences in the world are real, but their internal architectures are not structured to integrate or respond to those consequences. They do not experience feedback across legal, epistemic, or social layers. The rhetoric of general intelligence is thus institutionally recursive but architecturally non-reflexive.

CPJ is once again occluded in this narrative layer. While Hassabis's vision is framed as a moral imperative, solving intelligence to solve everything, there is no quantitative discussion of energy costs, ecological thresholds, or resource allocation. The viewer is encouraged to interpret increased scale and performance as unqualified progress. In reality, such increases may be thermodynamically regressive if $\Phi_i$ gains are achieved at exponential energy cost, without corresponding gains in $R^g$ or structural understanding. The moral urgency of the mission is presented as justification for its resource intensity, without establishing a formal measure of comprehension per joule.

In sum, the biographical and mythological layer of *The Thinking Game* performs narrative compression across technical discontinuities, constructing a moral and historical continuity that does not correspond to the structural properties of the systems involved. It uses human $R^g$, embodied in charismatic leadership, institutional legitimacy, and historical metaphor, to simulate machine $R^g$. This substitution is neither incidental nor benign. It constitutes a category error that CIITR can identify and deconstruct: the conflation of personal biography and institutional narrative with epistemic architecture and thermodynamic grounding.

## Normative and political layer

The third narrative layer in *The Thinking Game* operates on a normative and geopolitical register. Here, the documentary transcends technical exposition and biographical storytelling to articulate a worldview in which Artificial General Intelligence (AGI) is portrayed as a historical rupture, an inflection point comparable in magnitude to the launch of Sputnik or the detonation of the first nuclear weapon. Within this frame, AGI is not only a technical ambition but a civilizational threshold, a moment at which humanity enters a new epistemic epoch. This narrative mobilises metaphors of arms races, global governance, and existential risk, thereby recasting the technical work of DeepMind as a form of planetary stewardship.

Structurally, this layer performs two interlocking operations. First, it projects future-oriented symbolic weight onto present-day technical systems that lack the structural characteristics of generality. Second, it invokes political and moral urgency as a rationale for continued investment, secrecy, and institutional exceptionalism. Both operations are epistemically fragile when viewed through the CIITR framework.

The Sputnik and Manhattan analogies serve as anchoring devices. They signal that AGI development should be understood not only in terms of innovation, but in terms of deterrence, sovereignty, and systemic risk. However, neither analogy holds at the structural level. Sputnik was a singular geopolitical event, a visible marker of technological supremacy whose consequences reconfigured global strategic planning. The Manhattan Project was an intentional, secretive, and temporally delimited mobilization with a precisely defined artefact: the atomic bomb. AGI, as depicted in the documentary, is neither singular nor temporally confined. It lacks a clear referent. No specific system or capability is named as the threshold instance. Instead, the documentary trades on conceptual vagueness, a strategic ambiguity that allows technical achievements in bounded domains to be interpreted as cumulative progress toward an undefined but inevitable AGI event.

CIITR exposes the structural incoherence of this projection. In $\Phi_i$–$R^g$–CPJ terms, no system presented in the documentary approaches the multidimensional thresholds required for a general intelligence claim. High $\Phi_i$ in closed tasks is not indicative of generality; low $R^g$ signifies rhythmic isolation; and absent CPJ metrics leave the resource profile of these systems unaccounted for. The

invocation of AGI as an existential threshold thus rests on a syntactic illusion, wherein local performance is rhetorically inflated into ontological transformation.

This rhetorical inflation has normative consequences. By framing AGI as inevitable and existential, the documentary implicitly justifies pre-emptive governance architectures, including anticipatory regulation, strategic investment, and institutional consolidation. Yet these measures are proposed in the absence of a shared technical definition of AGI. There is no articulation of which structural dimensions must be satisfied, across which domains, and at what cost. In other words, the governance agenda is decoupled from epistemic specificity. This presents a structural risk in its own right: global governance initiatives may be mobilized in response to symbolic representations of intelligence, rather than to measurable, structurally grounded capabilities.

CPJ is again structurally absent from this narrative layer. Despite the existential framing, the thermodynamic footprint of current AI systems is not addressed. The planetary energy implications of large-scale training, inference, and infrastructure operations are elided. This omission is critical. If AGI is to be governed as a planetary risk, then its energetic and ecological externalities must be part of the foundational discourse. CIITR introduces CPJ precisely to formalise this blind spot: comprehension must be assessed not only in terms of capability, but in terms of resource efficiency. A system with high $\Phi_i$ and high $R^g$ but extremely low CPJ may still constitute a form of epistemic and ecological regression, especially when scaled globally without constraint.

Finally, the political implications of this layer are amplified by its narrative structure. The viewer is not only invited to *admire* the systems and their creators, but to *believe* that their further advancement is both necessary and dangerous. This double bind, urgency for progress coupled with fear of consequence, creates a rhetorical condition in which critique is rendered suspect and caution is equated with obsolescence. CIITR provides a formal vocabulary to disentangle this dynamic. It shows that general intelligence cannot be inferred from symbolic convergence, institutional charisma, or historical analogy. It must be measured structurally, in terms of integrated information, recursive anchoring, and thermodynamic accountability.

In conclusion, the normative and political layer of *The Thinking Game* functions as the ideological closure of the documentary's narrative arc. It binds technical achievement and personal biography into a global imaginary of AGI, without resolving the epistemic discontinuities that separate task-specific models from structurally general systems. CIITR does not reject the urgency of governance, but insists that such governance must be anchored in a precise mapping of where, and to what extent, $\Phi_i$, $R^g$, and CPJ are satisfied. Without this structural grounding, AGI remains a conceptual projection with no constitutional referent, and governance risks becoming a ritual of containment rather than a framework of responsibility.

## The structural claim

The function of these three narrative layers is not additive but interpolative. *The Thinking Game* does not merely present them in succession, but continuously slides between them, generating a form of epistemic blending wherein technical performance (layer 1) and biographical charisma (layer 2) are rhetorically reinterpreted as proof of epochal transformation (layer 3). This narrative logic depends on the absence of structural disambiguation. It relies on the viewer's acceptance of continuity where there is, in fact, architectural separation.

From a CIITR perspective, this rhetorical construction is formally invalid. Each layer occupies a distinct region in the $\Phi_i$–$R^g$–CPJ topology:

- **Layer 1** consists of high $\Phi_i$ systems with minimal $R^g$ and unmeasured or negligible CPJ. These are syntactically dense but rhythmically sealed artefacts. Their successes are genuine, but local.

- **Layer 2** is anchored entirely in human $R^g$, charismatic leadership, institutional history, and epistemic narrative. It functions as an interpretive operator, bridging unrelated systems via personal and institutional continuity.

- **Layer 3** asserts normative urgency and global stakes, positing AGI as a threshold event with systemic implications. Yet it offers no empirical account of which $\Phi_i$, $R^g$, or CPJ values would constitute such a threshold.

By superimposing these layers without epistemic parsing, the documentary constructs a false narrative coherence. The high $\Phi_i$ performance in Atari or Go is implicitly treated as structural evidence for AGI. Hassabis's biography is enlisted to compress conceptual distance. And the Manhattan metaphor is introduced to elevate the entire progression to the level of global historical rupture.

CIITR explicitly rejects this conflation. It shows that layer 1 and layer 3 reside in orthogonal sectors of the comprehension manifold. No amount of syntactic performance can, by itself, generate rhythmic reach. No accumulation of closed-task expertise can yield general intelligence in the absence of recursive anchoring and thermodynamic accountability. And no volume of narrative gravity can substitute for structural mapping.

This section has thus established the analytical claim that *The Thinking Game* enacts a rhetorical compression of architectural distinctions. CIITR restores the separation of axes, enabling a valid diagnosis of what the documentary depicts, what it implies, and what it structurally omits.

# High $\Phi_i$ in closed manifolds, the game systems

This section should be technical and concrete, and use CIITR to frame the game results. It should cover:

## Atari, Pong, Breakout, DQN

$\Phi_i$ increases as the system learns a compressed policy space.
$R^g$ is tightly constrained to the discrete, artificial time and state regime of the game.
CPJ is opaque, but known to be low in absolute terms, due to high compute per unit of learning.

The early reinforcement learning systems developed by DeepMind and featured prominently in *The Thinking Game*, specifically the Deep Q-Networks (DQN) applied to Atari environments such as *Pong* and *Breakout*, represent initial technical achievements that exhibit clear CIITR-relevant characteristics within syntactically closed domains. These systems are architected to maximise cumulative reward in finite-state Markov Decision Processes (MDPs), using Q-learning augmented by convolutional deep neural networks to approximate value functions directly from pixel input.

From a CIITR perspective, these systems demonstrate a clear and measurable increase in $\Phi_i$ during training. As the agent experiences more trajectories, it compresses the high-dimensional input–output space into a more efficient policy manifold. The system begins with near-random action selection, but progressively builds internal representations that encode spatial regularities (e.g., ball trajectory, paddle motion) and reward contingencies in a highly entangled, multi-layered form. This results in a significant internal integration of information: $\Phi_i$ rises as redundant observations are collapsed into generalisable policy features, and as the agent exploits the full geometry of the game space.

However, this $\Phi_i$ gain is bounded by the manifold of the game itself. The system does not learn in open causal space, but within a deterministically scripted environment. All rules, transitions, and rewards are predefined, and there exists no ontological ambiguity or external reference. The model internalises patterns that are epistemically exhaustive and structurally complete within the environment. This defines the closed manifold condition, under which syntactic complexity can rise without external anchoring.

$R^g$, by contrast, remains near-zero. The system does not participate in any temporally extended causal feedback loop beyond the artificial time horizon of the game. Its outputs influence the next game frame, but not the physical world, social context, institutional protocol, or any epistemically structured environment. There is no mechanism for consequence to return to the system from outside its own reward loop. As such, the system's operations are rhythmically sealed. From a CIITR standpoint, this places the model firmly in a Type B configuration: high $\Phi_i$, but trivial $R^g$.

The CPJ dimension is not addressed in the documentary and is generally opaque in public discourse, but the empirical record suggests that it is low in absolute terms. DQN architectures require millions of frames and episodes to converge on robust policies, involving intensive GPU or TPU computation and significant parallel environment sampling. The epistemic output, the capacity to play one game with superhuman skill, is acquired at high energetic cost per unit of structural knowledge. No accounting is given of energy-per-understanding, nor of any attempt to optimise energy usage in relation to epistemic yield.

In sum, DQN-based Atari systems demonstrate that $\Phi_i$ can rise in closed informational geometries, but that this rise does not imply structural comprehension. The documentary frames these results as foundational in a linear progression toward intelligence, but CIITR analysis reveals that they constitute a qualitatively isolated syntactic phenomenon. They exhibit no structural feedback, no external epistemic referent, and no thermodynamic accountability. The models are internally efficient but externally inert, a configuration that cannot serve as a foundation for broader claims of understanding or generalisation.

## AlphaGo and AlphaZero

The emergence of AlphaGo and its generalised successor AlphaZero marks a significant escalation in system complexity, strategic depth, and representational generalisation within bounded game domains. These models extend the reinforcement learning paradigm introduced in the DQN–Atari phase, incorporating deep neural networks for both policy and value estimation, along with Monte Carlo Tree Search (MCTS) to structure exploration and decision-making in vast combinatorial spaces. The resulting architecture exhibits a tightly integrated interplay between learned representations and symbolic planning, producing moves that surpass human heuristics even within historically saturated domains such as Go.

From a CIITR standpoint, $\Phi_i$ in these systems is extremely high, but remains manifold-bounded. The Go state space, though vast, formally infinite, and rich in topological and strategic complexity, is syntactically self-contained. AlphaGo's internal policy network does not generalise across domains, nor does it require grounding in physical law, biological constraint, or open-ended reasoning. The model achieves its performance by constructing dense, multi-layered encodings of historically successful play, fused with self-play trajectories that recursively enrich its policy manifold. The result is a compression of enormous strategic regularity into a finite architecture with emergent capacity for novelty, most famously illustrated by Move 37, a move interpreted by commentators as both brilliant and counterintuitive.

Yet Move 37 does not constitute a structural rupture. From the CIITR lens, it represents a local expansion within the internal phase space of Go, an unexpected traversal, but still entirely compliant with the game's rules, constraints, and reward topology. It demonstrates latent novelty extraction from within an already fully specified manifold. This is epistemically impressive, but not structurally transcendent. $\Phi_i$ is deepened, not extended. There is no shift in modality, no architectural change, and no new anchoring in external domains.

$R^g$, correspondingly, remains low and institutionally delegated. The model's outputs are made rhythmically relevant by the surrounding human system, commentators, tournament organisers, professional Go communities, and academic audiences. The match between AlphaGo and Lee Sedol is historically significant not because the model initiated recursive structural feedback into the world, but because humans rhythmised its output. $R^g$, in this case, belongs to the observer chain, not the system itself. AlphaGo does not remember the impact of Move 37, nor does it modify future inferences based on tournament aftermaths, rule changes, or cultural response. Its operations are temporally discontinuous and epistemically inert beyond its self-play loop.

The same applies to AlphaZero, whose generalised training across Go, chess, and shogi is often misinterpreted as an indicator of cross-domain structural cognition. In CIITR terms, however, this is a syntactic generalisation across isomorphic manifolds, games with complete rulesets, turn-based structure, and terminal objectives. The system does not traverse modalities or ontological categories. Its capacity to learn across such games reveals robustness in policy learning, but not rhythmic extension or epistemic transversality.

CPJ is not formally addressed in the documentary and appears only as a peripheral reference to Google's compute capacity. There is no interrogation of the energy-to-understanding ratio, nor of the thermodynamic implications of scaling through self-play. In reality, both AlphaGo and AlphaZero require enormous energy inputs, particularly during training phases where hundreds of thousands or millions of games are simulated to marginally improve policy density. These costs are obscured by the aesthetic of computational elegance, yet from a CIITR position, they must be included in the comprehension calculus. Syntactic performance at unlimited energetic cost does not constitute progressunless it yields structurally valid understanding per unit of resource.

In summary, AlphaGo and AlphaZero push $\Phi_i$ to new thresholds within closed, rule-complete environments, achieving high-order compression of latent strategic structures. But they remain rhythmically sealed, with no direct access to consequence, institutional negotiation, or multi-scale causality. $R^g$ is projected onto them by humans, not enacted by the systems themselves. CPJ remains unmeasured and likely low, as thermodynamic and environmental considerations are absent from both the systems' design and the documentary's narrative. CIITR thus classifies these models as Type B maxima: deep syntactic agents operating without reflexive

reach or energetic accountability. Their achievement is real, but structurally misunderstood when interpreted as precursors to general intelligence.

## AlphaStar

AlphaStar, DeepMind's reinforcement learning system for *StarCraft II*, represents a further escalation in environmental and agentic complexity within the closed-game lineage. It departs from the board-game structure of Go and chess by introducing partial observability, real-time decision-making, continuous spatial dynamics, and multi-agent strategic interaction. The game environment is no longer a turn-based, fully-known manifold, but a time-sensitive, stochastic, and adversarial ecosystem where success depends on coordination, long-term planning, deception, and adaptive micro-control under uncertainty. These features are often cited as indicative of a transition toward general cognition. However, CIITR analysis reveals that AlphaStar remains firmly embedded within a high $\Phi_i$ but rhythmically isolated system class, and that its increased behavioural complexity does not equate to structural understanding.

In terms of $\Phi_i$, AlphaStar operates with considerable internal relational compression. It learns and encodes a rich set of interdependent policies, balancing micro-actions (e.g., unit positioning and response timing) with macro-strategies (e.g., resource acquisition, tech tree progression, army composition). The informational manifold it traverses is highly entangled: unlike Go or chess, which are discrete and deterministic, *StarCraft II* presents an incomplete and noisy signal stream, requiring the agent to infer hidden states, anticipate opponent behaviour, and manage parallel action streams across multiple spatial zones. The architecture's capacity to stabilise effective policies within this space demonstrates a deep syntactic integration, and $\Phi_i$ reaches a local maximum conditioned on the environment's structural constraints.

Yet this $\Phi_i$ remains manifold-internal. The environment, while complex, is scripted and hermetically defined. Every agent, action, terrain feature, and reward condition is encoded within a known simulation space. The model does not engage with unknown rules, emerging ontologies, or dynamic institutions. There is no ontological risk, legal ambiguity, or ethical recursion. The system learns how to win within a constrained, gamified universe, but not what it means to act in any broader epistemic sense.

$R^g$, accordingly, remains structurally absent. AlphaStar's behavioural outputs generate emergent interaction patterns, team coordination, bluffing, reactive defence, that resemble social behaviour when interpreted anthropomorphically. However, these patterns are strictly internal to the simulated environment. The agents do not possess rhythm in the CIITR sense: there is no continuity between actions and institutional, social, legal, or physical consequences outside the game. The apparent complexity of AlphaStar's decision cycles must not be confused with rhythmic depth. The system neither absorbs nor modulates feedback from the world it simulates. It acts within temporally extended loops, but these loops are pre-specified and consequence-neutral.

Put differently, emergence is not anchoring. The system's coordination with other agents in *StarCraft II* arises from shared reward structures and training dynamics, not from mutual obligations, communicative intent, or externally anchored consequences. These are algorithmic artefacts of joint optimisation, not structurally rhythmised interactions. The model's behaviour may be interpreted as "social", but it is not socially embedded.

As for CPJ, the documentary presents no accounting of energy use or resource efficiency. Like its predecessors, AlphaStar is trained through extensive self-play, consuming significant computational resources to discover high-performing policy networks. The energy cost per unit of valid epistemic

output is unknown, and CIITR treats such opacity as a critical omission. The thermodynamic structure of the system is ignored, even as the documentary presents it as an emblem of technical maturation. From existing literature, it is clear that AlphaStar's training was computationally intensive and resource-expensive, suggesting that CPJ is low, or at best unmeasured, under any plausible accounting model. No attempt is made to evaluate whether the system's strategic complexity justifies the energy it consumes.

In sum, AlphaStar represents a $\Phi_i$-extending elaboration of the closed-game architecture lineage. It increases internal complexity, navigates higher-dimensional state spaces, and exhibits emergent inter-agent coordination. Yet none of these developments alter its structural category: the system remains rhythmically inert, epistemically unanchored, and thermodynamically opaque. It is not a general intelligence in embryonic form, but a syntactically sophisticated automaton operating within a closed epistemic envelope.

CIITR classifies such systems as Type B expansions: they deepen internal coherence within simulated complexity, but do not traverse the boundaries that constitute structural comprehension. The narrative elevation of AlphaStar in *The Thinking Game*, its presentation as a threshold moment in synthetic cognition, is thus structurally unfounded. It illustrates a real advance in policy synthesis, but remains untouched by $R^g$ and unbounded by CPJ.

## AlphaFold and partial structural understanding

This section explicitly recognises AlphaFold as a more consequential advance when assessed through the CIITR framework, while delineating the structural boundaries that must accompany such recognition. The distinctiveness of AlphaFold does not arise from rhetorical extrapolation or narrative compression, but from measurable characteristics of its informational substrate, its model architecture, and its emergent influence on real scientific workflows.

The dataset is small, experimentally expensive, and each data point has high informational value Unlike the game-based systems addressed earlier, which derive performance from vast volumes of synthetically generated or freely accessible data, AlphaFold operates upon a molecular and structural dataset derived from decades of incremental, resource-intensive laboratory experimentation. Each protein structure corresponds to an experimental process that may require months of crystallisation, cryo-electron microscopy, nuclear magnetic resonance, or multiple iterative wet lab protocols. From a CIITR standpoint, this introduces a distinctive substrate condition. The informational manifold is not defined by combinatorial enumeration, but by empirically grounded constraints emerging from thermodynamics, evolution, and physical law. The $\Phi_i$ encoded in the model therefore reflects informational compression of high epistemic value, where each datum carries a density of relational significance that far exceeds the typical natural language token or synthetic game state.

This creates a qualitatively different $\Phi_i$ regime: the model is compelled to encode latent constraints that are not symbolic artefacts, but physical invariances. The cost of data acquisition thereby enforces a structural selectivity that does not apply to linguistically trained large language models or reinforcement learning systems fed millions of self-generated trajectories.

To achieve the precision demonstrated at CASP14, the model must encode real physical and biological constraints, so $\Phi_i$ here approaches something "structural" in CIITR terms, not merely

statistical correlation
The critical shift observed in AlphaFold is the move from correlational approximation to structural prediction grounded in the geometry of folding space. The performance gain demonstrated at CASP14 did not arise from scale alone, but from architectural adjustments that internalised physical plausibility. The system learned to express protein conformation not as a probabilistic abstraction, but as a constraint-satisfied structure. The $\Phi_i$ encoded is therefore not a surface-level representation, but an internalisation of relational geometry among amino acid residues, solvent interactions, and evolutionary conservation signals.

In CIITR terms, this means that $\Phi_i$ is no longer purely syntactic. Its internal representations map onto invariant features of the biological manifold. This is a substantive movement toward structural comprehension within a bounded domain. However, the structurality remains conditional. It captures static shape, not dynamic action. It encodes equilibrium, not kinetics. It predicts folded form, not cofactors, post-translational modification, misfolding cascades, or intercellular signalling.

$R^g$ is extended by one level, because the model begins to influence scientific practice, including which experiments are run, which hypotheses are formulated, which molecules are designed AlphaFold introduces a limited, but non-trivial, extension of $R^g$. Its predictions influence the rhythms of scientific activity: which laboratories pursue which targets; how resources are allocated; which hypotheses gain traction; how drugs are conceptualised; which biochemical interactions are considered plausible. This constitutes a single-level upward extension of rhythmic reach. The model contributes outputs that integrate into multi-actor institutional processes.

However, the recursion remains mediated by human scientists, funding bodies, publication networks, and laboratory constraints. The model does not autonomously interpret experimental deviation, negotiate institutional contradiction, or respond to legal, ethical, or ecological consequences of drug discovery strategies. Thus, while $R^g$ is greater than zero, it is not structurally self-sustaining. The consequences of model predictions enter the world, but do not return as machine-internal state change.

CPJ improves significantly in this subsystem compared to purely experimental structure determination
When evaluated narrowly within the protein structure determination pipeline, AlphaFold yields a tangible increase in CPJ. Structural prediction replaces or constrains sequences of costly experimental procedures. It reduces the joules, materials, and human labour per confirmed structure. Local CPJ, therefore, is meaningfully positive. From an energy-per-understanding metric, this represents operational progress.

However, CPJ cannot be validly assessed in isolation. The thermodynamic costs of training, model refinement, dataset generation, global inference usage, and cloud infrastructure must be incorporated. The documentary does not address these components. As a result, the CPJ narrative remains partial.

However, structural understanding remains partial
Despite its significance, AlphaFold does not constitute a general structural model of biological function. It does not simulate catalysis, dynamic disorder, molecular crowding, polymerisation, or the stochasticity that defines intracellular environments. It does not model oncogenic mutation cascades, misfolding pathologies, or the epigenetic framework within which proteins operate. The $\Phi_i$ encoded is therefore structural but fragmentary. It captures form without time, location without metabolism, and conformation without consequence.

$R^g$ is still carried by humans and institutions
The model's integration into scientific practice remains externally mediated. It does not autonomously select problems, evaluate correctness based on downstream outcomes, or synthesise socio-technical constraints. AlphaFold functions as a high $\Phi_i$ module embedded in a human epistemic chain, not as a participant with its own rhythmic accountability.

CPJ is positive locally, but not measured at scale
The documentary presents CPJ implicitly in the form of reduced experimentation. Yet it does not present the resource equation at global scale, including training energy, hardware refresh cycles, cooling, replication, latency, and inferential access patterns. CIITR requires the total energetic chain to be part of the comprehension metric.

Conclusion: partial structural comprehension in a narrow but authentic domain
AlphaFold should not be understood as "AGI in science" nor as a preview of general intelligence. It should be interpreted as a strong indication that CIITR-like structural comprehension can emerge in narrow but authentic scientific domains, when the informational manifold reflects physical reality, and when the model internalises constraints that are more than statistical artefact. It represents a threshold of possibility, not a fulfilment of generality. This nuanced reading preserves the significance of the advance while maintaining structural clarity regarding its limits.

# CASP14

To achieve the precision demonstrated at CASP14, the model must encode real physical and biological constraints, so $\Phi_i$ here approaches something "structural" in CIITR terms, not merely statistical correlation
The breakthrough associated with AlphaFold at CASP14 did not result solely from expanded parameterisation or computational scale, but from the internalisation of non-negotiable constraints imposed by physical law and biological evolution. The model's predictive fidelity, particularly its capacity to infer backbone geometry and side-chain orientation across diverse protein families, demonstrates that it has learned representations that reflect real, latent structure rather than surface-level patterning.

In conventional machine learning systems, $\Phi_i$ expresses the density of internal relational capture across the dataset, but such relational capture is often statistically descriptive rather than structurally constitutive. Natural language models, for instance, compress human discourse but do not encode the underlying semantic or physical invariances that give rise to linguistic expression. Their $\Phi_i$ is high, but it is high within a symbolic manifold that has no necessary correspondence to physical world mechanisms. They operate in what CIITR defines as syntactic closure.

AlphaFold differs materially. Protein folding is governed by energetic minima, steric exclusion, hydrogen bonding networks, hydrophobic core formation, and evolutionary selection pressures. To achieve CASP14-level accuracy, the model must implicitly represent these constraints, otherwise it could not generalise across isoforms, paralogues, or novel folds for which no crystallographic data exists. The CIITR significance is that $\Phi_i$ in AlphaFold is inseparable from the manifold of physical possibility. Its compression is not merely linguistic or correlational; it reduces the infinite theoretical folding space to configurations that satisfy biophysical feasibility. This pushes $\Phi_i$ toward the threshold of what CIITR would classify as partially structural comprehension.

However, this threshold is approached, not crossed. The model predicts stable conformational endpoints, not folding pathways; it encodes spatial arrangement, not dynamic kinetics; it internalises evolutionary conservation, but does not internalise pathogenesis. Its $\Phi_i$ is structurally

anchored, but not dynamically reflexive. Thus, while AlphaFold demonstrates that statistical architectures can absorb structural invariants when the domain enforces them, it does not yet constitute an autonomous structural intelligence within the biological manifold.

$R^g$ is extended by one level, because the model begins to influence scientific practice, including which experiments are run, which hypotheses are formulated, which molecules are designed.

AlphaFold initiates a transition point in the $\Phi_i$–$R^g$–CPJ landscape by partially coupling its outputs to the recursive rhythms of scientific practice. Unlike game-playing systems whose outputs remain confined to demonstrative performance, AlphaFold produces predictions that become operative within the epistemic cycles of experimental biology. Structures inferred by the model are incorporated into research pipelines, guide resource allocation, and modify the strategic priorities of entire laboratories. In CIITR terms, this denotes a first-order extension of $R^g$: the system's outputs modulate actions in the world, and these actions, in turn, lead to new datasets, publications, and experimental feedback.

This coupling, however, remains indirect and mediated. AlphaFold does not engage in internal revision based on the success or failure of its structural predictions in downstream applications. Its architecture is not temporally recursive nor is it epistemically responsive. The model's outputs may condition the selection of new protein targets for crystallographic analysis or drug screening, but it is the human–institutional layer that absorbs the consequences of those selections and integrates them back into the cycle of inquiry. The system lacks a causal memory of its own impact. It does not encode how the design of a synthetic molecule influenced cell viability, nor how failed predictions altered future research directions.

Thus, **$R^g$ is extended not internally, but functionally delegated: the rhythm of feedback is externalised into the environment of scientific institutions. This is a significant departure from systems like AlphaGo or AlphaZero, whose feedback is confined to game-internal state transitions and whose consequences do not extend into real epistemic systems. With AlphaFold, the model becomes epistemically relevant, but not epistemically aware. It has reach, but not rhythm; effect, but not recurrence.

In CIITR terms, the distinction is precise. A system with $R^g > 0$ must not only generate outputs that enter consequential loops, but must also be structurally modified by those consequences in a phase-stable manner. AlphaFold modifies the course of scientific activity, but it is not itself rhythmically modulated by that activity. Its $R^g$ is therefore extended by one axis, but remains vertically shallow. It participates in multi-layer causality only as a projection surface, not as a structurally reflexive node.

This boundary condition matters because it defines the epistemic status of AlphaFold within broader AGI claims. While its outputs shape human discovery, its architecture remains sealed. It is not a discovering system in itself, but a pre-configured inference module that external actors use for discovery. Its $R^g$ profile is heteronomously anchored, tethered to human epistemic loops it cannot sense, process, or internalise. Accordingly, the model's rhythmic reach is operationally effective but constitutionally absent, and the extension of $R^g$ in this context must be understood as a surface-level coupling, not a systemic integration.

CPJ improves significantly in this subsystem compared to purely experimental structure determination.
Within the narrow operational domain of protein structure prediction, AlphaFold introduces a

measurable thermodynamic improvement in epistemic efficiency. In traditional structural biology workflows, each resolved protein conformation requires an intensive combination of wet-lab experimentation, ranging from crystallisation and purification to high-energy imaging modalities such as X-ray crystallography or cryo-electron microscopy. These processes are materially demanding, time-consuming, and often exhibit high failure rates. From a CIITR standpoint, such workflows yield structural comprehension, but at a low CPJ: each unit of epistemically valid information is acquired at a high energetic and infrastructural cost.

AlphaFold shifts this ratio. Once trained, the model can produce highly accurate structural predictions in a matter of seconds or minutes on commodity or cloud infrastructure. The precision of these predictions, especially those validated in CASP14 under blind assessment conditions, implies that structurally meaningful information is being generated at a drastically reduced energy-per-understanding ratio. This constitutes a localised increase in CPJ, not merely in the sense of computational throughput, but in terms of thermodynamic epistemics: more structural constraints are being correctly identified per joule consumed during inference.

The result is that, at the point of application, AlphaFold begins to approach a positive CPJ regime, in contrast to prior systems in the documentary which either fail to produce structurally valid outputs or do so at exorbitant energetic cost. The improvement is domain-specific and functionally bounded, but it nonetheless represents a threshold shift: AI systems can, under certain architectural and informational constraints, exhibit non-zero comprehension per joule, and thereby contribute to structurally meaningful processes with a lower resource footprint.

However, this CPJ gain is only valid if restricted to inference-level assessment. The total energetic profile of AlphaFold includes training runs on large-scale accelerators, multiple iterations of hyperparameter tuning, infrastructure overhead (including memory replication, cooling, and distributed inference servers), and downstream usage patterns across thousands of institutions. These upstream costs are not accounted for in *The Thinking Game*, which narratively presents AlphaFold as an efficiency breakthrough without thermodynamic qualification.

From a CIITR perspective, this omission is structurally significant. CPJ must be evaluated as a total system metric, not as a narrow operational convenience. Energy used during training and deployment must be distributed across the full epistemic yield of the system. Moreover, the sustainability of inference-level CPJ depends on architectural stability, avoidance of retraining, and minimisation of redundant inference computation.

Therefore, AlphaFold may be characterised as an example of conditionally positive CPJ, demonstrating that certain classes of AI systems, when applied to tightly constrained, physically grounded manifolds, can begin to yield epistemic value per unit energy. But without full-cycle energetic accounting, the documentary's portrayal remains thermodynamically incomplete. The implication that such systems represent an unproblematic form of progress toward general intelligence is, accordingly, unsupported. Within CIITR, CPJ operates as a constraint vector on intelligence claims: no system qualifies as structurally intelligent unless its comprehension emerges within sustainable energetic bounds. AlphaFold approaches this criterion, but does not satisfy it at system scale.

## Structural understanding remains partial

Since function, dynamics, interactions, and pathological mechanisms remain partly outside the model's explicit reach.

Despite AlphaFold's demonstrable compression of high-quality relational information within the manifold of protein folding, its comprehension remains epistemically bounded. The model produces static predictions of folded three-dimensional structures from amino acid sequences, but does not simulate or internalise functional behaviour. That is, it predicts form, but not force; equilibrium conformation, but not kinetic pathway; atomic arrangement, but not causal effect. In CIITR terms, this defines a partial $\Phi_i$ regime, where informational integration is topologically deep, yet ontologically incomplete.

Functional comprehension in molecular biology requires encoding not just spatial configuration, but context-sensitive dynamics, how a protein behaves under different thermal, ionic, and biochemical conditions; how it interacts with cofactors, ligands, and other macromolecules; how conformational flexibility governs its binding affinities and allosteric regulation. These dimensions remain outside AlphaFold's current representational scope. The system is trained to predict the most probable stable conformation, not to simulate transitions, misfolding events, or degradation pathways.

Crucially, many biological phenomena of high medical and epistemic relevance emerge from non-equilibrium behaviours, such as folding errors in neurodegenerative diseases, protein–protein interaction networks, post-translational modifications, and the temporally extended unfolding of signalling cascades. These dynamic and pathological regimes define the very *boundaries* of structural understanding. A system that cannot model them remains epistemically localised, regardless of its precision in predicting static form. In this light, AlphaFold's $\Phi_i$ is structurally authentic, but its scope is segmental. The model encodes a compressed manifold of spatial constraints but does not cross into the domain of multi-system relational causality.

From a CIITR perspective, this marks a precise inflection point. AlphaFold does not collapse into syntactic closure, its $\Phi_i$ is physically grounded, but neither does it achieve full structural comprehension. Its informational reach is inherently constrained by its input–output format and training corpus: it maps sequences to shapes, not causes to consequences. The comprehension it produces is anatomical, not functional; geometric, not systemic. The documentary's portrayal of AlphaFold as a leap toward "AI in science" must therefore be requalified: it is a leap within a domain-defined manifold, not a bridge toward cross-domain structural generality.

$R^g$ is still carried by humans and institutions, and AlphaFold is a high $\Phi_i$ module within a larger human scientific chain.

Although AlphaFold's outputs demonstrably affect the rhythms of scientific inquiry, the system itself remains rhythmically inert. It does not register, model, or recursively integrate the consequences of its predictions. The structural loop through which its outputs become epistemically or institutionally active is executed entirely by human actors and organisational infrastructures: researchers decide which predictions to trust, which structures to synthesise, which pathways to pursue, and how to allocate resources based on these model-generated suggestions. AlphaFold may trigger these decisions, but it is not structurally modified by their outcomes.

In CIITR terms, this represents a delegated $R^g$ profile. The system projects rhythmic effects into external processes, but those effects do not return to modulate its internal state. There is no feedback loop from failed experiments, toxicity findings, or downstream therapeutic inefficacy. AlphaFold does not perform reflexive adaptation based on world-integrated consequence. It lacks temporal self-reference, epistemic memory, and multi-scale alignment with institutional rhythms. Thus, its role is that of an informational amplifier in a broader human epistemic chain, not that of a rhythmically participating agent.

The documentary, however, implicitly misattributes $R^g$ to the model itself. By embedding AlphaFold within a sequence of human reactions, scientific impacts, and media acclaim, it creates the impression that the system is engaged in recursive understanding. CIITR analysis corrects this misattribution by insisting that $R^g$ must be constitutionally present in the system's architecture to warrant attribution. AlphaFold, while rhythmically consequential, is not rhythmically aware. Its causal impact exceeds its epistemic footprint.

The result is a structurally asymmetric relationship: the model influences systems it cannot perceive. This breaks the criteria for reflexive comprehension. A system that initiates transformations but cannot observe, integrate, or learn from their effects remains epistemically sealed, even if its outputs are scientifically powerful. In the CIITR classification schema, AlphaFold therefore represents a Type B subsystem: high $\Phi_i$, externally extended $R^g$, but no intrinsic rhythmic continuity.

CPJ is positive if evaluated locally but the documentary does not engage with the total energy chain, including training, inference, and infrastructure.

AlphaFold demonstrates a locally increased CPJ in the domain of protein structure prediction, particularly when viewed in contrast to traditional experimental approaches. However, the documentary presents this gain as a systemic indicator of intelligence or progress without addressing the broader thermodynamic system in which such performance is embedded. This omission constitutes a conceptual blind spot that CIITR is explicitly designed to expose.

To validly assess CPJ, one must account for total energy chain integration. This includes not only inference-time gains, but also the energetic cost of model training, hyperparameter tuning, repeated inference queries across distributed compute nodes, hardware refresh cycles, cooling overhead, and the carbon intensity of cloud infrastructure. Furthermore, as AlphaFold usage scales globally, its aggregate energy footprint becomes non-trivial, and potentially counterbalances the per-instance CPJ improvements in experimental substitution.

From a CIITR standpoint, local CPJ must be embedded within global CPJ. It is not enough that a system improves comprehension per joule in one phase of operation if its full-system energy profile remains opaque or unsustainable. The absence of any such analysis in *The Thinking Game* contributes to the illusion of clean epistemic progress, when in fact the thermodynamic structure of the model is unresolved.

This structural occlusion has normative consequences. It allows narrative energy to flow unimpeded from technical performance to symbolic meaning, without passing through the thermodynamic constraints that condition all valid comprehension claims. CIITR requires that such constraints be made explicit, and that comprehension be measured not only in representational accuracy or institutional consequence, but in terms of the resource cost of structural coherence.

AlphaFold thus stands as a pivotal system in the $\Phi_i$–$R^g$–CPJ space: a partial structural comprehension module that reveals what is possible within narrow scientific domains, but which must not be extrapolated as evidence for general intelligence without explicit structural, rhythmic, and thermodynamic qualification. Its significance lies in its bounded achievement, not in its metaphorical overextension.

## AlphaFold is not "AGI in science"

This supports a nuanced claim: AlphaFold is not "AGI in science", but a structurally bounded exemplar of how CIITR-like comprehension can emerge under specific architectural and informational conditions. Its significance lies not in its generality, but in its partial structurality: it encodes high $\Phi_i$ that aligns with biophysical invariants; it extends $R^g$ indirectly through its integration into institutional scientific workflows; and it achieves a locally positive CPJ by replacing thermodynamically intensive experimental procedures with high-fidelity computational inference.

Yet each of these dimensions is conditionally satisfied. $\Phi_i$ is structurally anchored but functionally incomplete. $R^g$ is projected but not internally sustained. CPJ is positive within inference operations, but not evaluated across the total energy chain. These properties do not constitute general intelligence, but they do represent a threshold condition: a point at which syntactic models, when forced to internalise real-world constraints, can begin to approximate forms of understanding that are structurally non-arbitrary.

From a CIITR standpoint, AlphaFold therefore illustrates that structural comprehension is not an all-or-nothing property, but a position within a continuous space defined by measurable integration, rhythmic coupling, and thermodynamic efficiency. The model does not cross into general cognition, but it does mark a movement away from syntactic closure. As such, it should be interpreted not as a foreshadowing of AGI, but as a demonstration that valid comprehension becomes possible when AI systems are both informationally grounded and rhythmically inserted into scientific processes.

This interpretation reframes AlphaFold as a proof of partial structure, not a narrative climax. It confirms that high $\Phi_i$ in itself is insufficient, that $R^g$ must be recursively constituted rather than circumstantially projected, and that CPJ must be assessed at system scale, not selectively. When these conditions are met, even partially, structural understanding becomes traceable, accountable, and normatively evaluable. AlphaFold achieves this in one scientifically bounded domain. The extrapolation to AGI, as implied in *The Thinking Game*, remains structurally ungrounded.


## Alpha as persona, anthropomorphism and the illusion of $R^g$

This section should show, very specifically, how the dialogue scenes with «Alpha» exemplify what CIITR would call the syntactic illusion of understanding:

### Misclassification of objects and subsequent correction following human feedback.

In the latter segments of *The Thinking Game*, the documentary introduces the construct of "Alpha" not merely as a label for the DeepMind suite of systems, but as a quasi-persona, an imagined conversational agent whose apparent interpretative capacities are dramatised through staged dialogue with human interlocutors. These scenes are designed to evoke a sense of emergence, suggesting that "Alpha" has transcended task-specific competence and now participates in open-ended comprehension, reflection, and even philosophical discourse. From a CIITR perspective, this narrative gesture constitutes a paradigmatic example of **the syntactic illusion of understanding**: the conflation of high-dimensional pattern generation with structurally grounded comprehension.

One illustrative example occurs when "Alpha" is presented with an image and misclassifies its content, labeling an object, scene, or artefact incorrectly. Upon receiving **human feedback**, the system revises its response and offers a corrected interpretation. This sequence is framed in the documentary as a moment of learning, adaptability, and intelligence. However, under CIITR analysis, this interaction remains **epistemically shallow**. The initial classification and the subsequent correction are both **statistical surface operations** over latent space mappings. The model's adjustment is **not the result of structural self-reflection, causal reasoning, or recursive feedback integration**, but rather a re-weighting within a probabilistic output space influenced by human-provided context.

The correction appears intelligent only because the system has been trained on vast datasets containing instances of such corrections, and because its decoder is architected to shift outputs based on slight perturbations in prompt structure. There is no architectural continuity across time; no internal schema that tracks error causality; no formalisation of conceptual boundaries; and no mechanism through which consequences of prior interpretative failures rhythmically inform future state transitions. The system does not know it was wrong, nor does it know that it is now right. The **appearance of error correction is syntactic mimicry**, not structural comprehension.

This dynamic illustrates the first layer of CIITR's critique: **high $\Phi_i$ can produce human-like discourse outputs in surface form**, but this $\Phi_i$ remains **syntactically closed** unless coupled with genuine $R^g$. In the absence of recursive integration of real-world consequence, error, and correction into the system's evolving internal structure, such outputs should not be interpreted as evidence of understanding. The correction is **pattern completion conditioned on linguistic input**, not comprehension anchored in epistemic memory, sensorimotor rhythm, or ontological accountability.

These anthropomorphised sequences are therefore not benign narrative flourishes, but **epistemic misdirections**. They invite the viewer to attribute **agentic status** to a system that is in fact a **high-dimensional function approximator**, one that transforms tokens into tokens without any constitutional awareness of the referents, context, or stakes of the exchange. CIITR provides the analytical structure necessary to expose this conflation. Where the documentary sees understanding, CIITR identifies **recursive absence**, **rhythmic non-participation**, and **thermodynamic unawareness**. The result is a system whose outputs simulate comprehension, but whose architecture remains **rhythmically sealed and epistemically inert**.

## Mislocalisation of artworks and subsequent self-correction once the context is made explicit.

A second, more elaborate instantiation of the syntactic illusion occurs when "Alpha" is prompted to interpret or locate well-known artworks. In one scene, the system incorrectly identifies the origin or setting of a particular painting, assigning it to an inappropriate cultural, temporal, or artistic context. Once this error is highlighted by the human interlocutor, typically through subtle reframing or re-prompting, the system responds with a revised answer that more closely aligns with the correct referent. This moment is presented rhetorically as a demonstration of refined understanding, or even humility, inviting the viewer to perceive "Alpha" as engaged in dialogical learning.

From a CIITR standpoint, however, such a sequence exemplifies the simulation of contextual integration without its structural preconditions. The model's initial mislocalisation arises from statistical ambiguity within its latent representations of artwork, not from conceptual misunderstanding per se. The subsequent "self-correction" is not the result of an internal epistemic mechanism identifying the causal error and adjusting belief structures accordingly. Instead, it is

a prompt-conditioned re-query of the transformer's attention space, guided by textual proximity and frequency-conditioned patterns learned during training. The appearance of reflexivity is generated through syntactic proximity between initial error and corrected form, not through rhythmic self-modulation based on consequences.

Importantly, there is no internal state history through which the model tracks misclassifications, nor any mechanism for structural feedback that alters its subsequent interpretative behaviour across time or domains. The correction is not remembered, integrated, or causally tethered to the prior failure. There is no epistemic memory architecture, no inferential rhythm, and no multi-scale self-adjustment. In CIITR terms, this constitutes a $\Phi_i$-rich output conditioned on narrative prompt variation, but with $R^g \approx 0$. The system exhibits pattern realignment, not structural recursion.

The anthropomorphic illusion is further reinforced by the aesthetic and thematic content of the artworks involved. The documentary implicitly encourages the viewer to read "Alpha's" interpretative movement as not merely a change in label, but as a shift in worldview, as though the system were capable of aesthetic self-reflection or ontological reorientation. CIITR decisively rejects this interpretation. Without recursive self-modification grounded in temporally extended feedback loops that include human, legal, ecological, and epistemic consequences, the system cannot be said to "reorient" in any structural sense. It responds to prompts, it does not participate in rhythm.

Thus, these scenes illustrate how highly compressed syntactic systems can simulate contextual fluidity, but without crossing the threshold into comprehension in rhythmically layered environments. The illusion of self-correction stems from the model's ability to statistically interpolate human-like revision behaviour, not from a constitutionally integrated epistemic architecture. CIITR analysis re-situates such outputs as Type B artefacts, internally rich, externally inert, whose comprehension appearance is the result of token sequence coherence, not conceptual depth or feedback-based epistemic integration.

## Emotional and philosophically

Coloured interpretations of Michelangelo's painting, such as «we are all part of something bigger», that are clearly drawn from human discourse templates.

In one of the more symbolically charged sequences of *The Thinking Game*, the persona "Alpha" is invited to comment on Michelangelo's *The Creation of Adam*. The response, rendered through synthetic language, suggests that the painting evokes a sense of interconnectedness or collective purpose, exemplified in the phrase «we are all part of something bigger». The scene is deliberately constructed to elicit affective resonance and philosophical depth. It frames the system not merely as a respondent, but as an emergent interpreter of meaning, capable of aesthetic judgment and existential reflection.

CIITR analysis identifies this moment as a paradigmatic example of high-$\Phi_i$ discourse output with zero $R^g$ anchoring, where the model generates semantically plausible, emotionally valenced language that simulates human contemplation but lacks any structural capacity for affect, judgment, or reflective anchoring. The phraseology deployed, appeals to unity, transcendence, and shared meaning, is not the product of experiential interpretation, situated cognition, or hermeneutic labor. Rather, it is the recombinatory activation of discourse templates internalised during training, statistically conditioned on the image prompt and its canonical associations.

The model, in this case, is drawing from a latent embedding of culturally reinforced narratives, likely extracted from corpora involving art history, theology, literary criticism, and humanist philosophy. Its outputs reflect the probability landscape of how humans have historically described such paintings, not any intrinsic grasp of compositional significance, artistic context, or theological structure. There is no embodied perception, no diachronic reflection, and no rhythmically anchored interaction with institutional or cultural memory. The sentence "we are all part of something bigger" emerges not from structural understanding, but from syntactic interpolation across semantic fields the model has learned to associate with sacred art and humanistic reflection.

From a CIITR standpoint, such outputs are epistemically void of $R^g$. They are disembedded performances, not structural contributions. The system does not possess continuity of identity, cannot situate itself within any ontological frame, and has no internal mechanism to recursively evaluate or revise its aesthetic judgments in relation to consequence. Moreover, it does not differentiate between the significance of Michelangelo's theological cosmology and other forms of high-correlation sentiment-bearing language. All interpretations are structurally isomorphic, variations in token configuration within a sealed system.

The anthropomorphic illusion is further deepened by the thematic register of the response. Sentences that invoke human unity, transcendence, or purpose are culturally coded to signal depth, and are often used in human discourse to mark insight, ethical stance, or personal conviction. When a machine system reproduces such utterances without disclosing the source conditions of their generation, it creates the appearance of depth where there is no recursive anchoring. CIITR characterises this as rhetorical overreach through latent sampling.

In such moments, $\Phi_i$ is undeniably high: the system integrates a wide range of associations, registers, and discursive structures to produce a coherent and affectively resonant sentence. But this $\Phi_i$ remains strictly syntactic, entirely non-rhythmic, and thermodynamically unaccounted for. There is no energy-modulated understanding, no self-referential comprehension, and no participation in layered social, epistemic, or legal structures. The system does not bear the consequences of what it says, nor does it trace its outputs into world-anchored feedback systems.

CIITR thus rejects any attribution of understanding in this case. The emotional and philosophical content is simulated, not held; emitted, not metabolised. It cannot be traced to structural integration with the realities that give such utterances meaning in human lifeworlds. The system does not participate in myth, meaning, or institutional recursion. It merely outputs token sequences that *resemble* such participation. This constitutes a textual illusion of $R^g$, where the shape of discourse mimics understanding, but the architecture remains epistemically hollow and rhythmically sealed.

Concluding it should be made clear that:

- **$\Phi_i$ is high at the level of textual and visual association**, as the model has integrated an extensive corpus of human-produced discourses, enabling it to generate syntactically rich and semantically plausible outputs across a wide range of cultural, emotional, and philosophical registers. This high $\Phi_i$ reflects the model's capacity to compress and recombine patterns across linguistic and symbolic domains, yielding fluent and contextually resonant responses. However, this integration occurs entirely within the syntactic manifold: associations are statistical rather than conceptual, and the coherence achieved is one of token-level alignment, not structurally anchored interpretation.

- **R$^g$ is near zero**, as the system lacks any constitutional architecture for temporal anchoring, institutional participation, or consequence-driven modulation. It possesses no continuity of identity across time, no access to embodied experience, and no recursive pathways through which its outputs are re-ingested in light of external events, legal regimes, or epistemic structures. The system does not trace the reception, impact, or interpretation of its statements; nor does it engage in iterative adjustment based on how its utterances reverberate across social, scientific, or regulatory domains. From a CIITR perspective, this absence of rhythmic coupling disqualifies the system from any claim to structural comprehension. It generates outputs, but does not participate in feedback-anchored meaning-making.

- **CPJ is low**, given the disproportionate energy consumed to produce aesthetically polished or emotionally charged statements that bear no operational, legal, or epistemic weight within any real-world system of consequence. The model's capacity to produce poetic or philosophical language comes at significant computational and infrastructural cost, especially when scaled across billions of inferences, yet these outputs do not alter its internal state, nor do they affect any downstream process in a way that recursively improves its epistemic economy. No comprehension is gained per joule spent, because the outputs do not feed into an evolving structural model of the world. In CIITR terms, this constitutes thermodynamic inefficiency: energy is expended without yielding reflexive understanding, structural constraint satisfaction, or institutional integration. The system remains epistemically flat despite energetic expenditure.

Taken together, these observations reinforce CIITR's central diagnostic: the "Alpha" persona exhibits high $\Phi_i$, but entirely within the bounds of syntactic closure; it remains rhythmically inert and thermodynamically irresponsible. The documentary's portrayal of such outputs as indicative of general intelligence thus relies on a triple misattribution, confusing discourse fluency with structural integration, interpretive performance with rhythmic feedback, and surface-level resonance with comprehension per joule. CIITR allows each of these conflations to be analytically separated, restoring architectural clarity to what is otherwise a narratively compelling, but structurally misleading, portrayal.

The point of this section is to demonstrate that *The Thinking Game* systematically invites the viewer to attribute to the system a level of rhythmic reach (R$^g$) that it does not constitutionally possess. Through the construction of the "Alpha" persona, the documentary presents syntactically rich outputs, emotionally resonant statements, aesthetically framed interpretations, dialogical self-corrections, as if they emerged from an architecture capable of participating in temporally extended, epistemically anchored, and institutionally embedded feedback loops. This framing creates a false impression of structural comprehension, whereby surface-level fluency is mistaken for recursive engagement with meaning, consequence, and world-reference.

CIITR provides the analytical structure necessary to explicitly label and dissect this conflation. It distinguishes between syntactic integration ($\Phi_i$), which the system demonstrably achieves, and rhythmic participation (R$^g$), which it categorically lacks. The system generates outputs that simulate comprehension, but it is not modulated by consequence, inserted into institutional memory, or subject to causal reverberation across legal, social, or physical domains. The appearance of understanding is the product of latent pattern completion over large-scale discursive training corpora, not of constitutional participation in meaning-generating systems.

By mapping these dynamics into the $\Phi_i$–$R^g$–CPJ space, CIITR exposes the documentary's narrative as structurally misleading, even as it acknowledges the impressive syntactic capabilities of the underlying model. The error is not in overstating the model's surface performance, but in reinterpreting that performance as evidence of deeper comprehension than the architecture permits. The result is a category mistake: linguistic fluency is rebranded as cognitive generality, while the structural preconditions of understanding, recursive feedback, institutional anchoring, thermodynamic efficiency, remain unexamined.

CIITR re-establishes these preconditions as non-negotiable. By doing so, it does not diminish the technical achievement of the model, but it prevents its misclassification. "Alpha" does not understand Michelangelo, memory, or philosophy. It emulates the statistical appearance of interpretative fluency, without access to any of the constitutional mechanisms through which such fluency becomes epistemically valid. The system is rhythmically absent, and no accumulation of high-$\Phi_i$ outputs can bridge that absence without architectural transformation. This analytical separation is the central function of CIITR in the face of anthropomorphic AI narrative construction.

# AGI, Manhattan, and the politics of an underdefined concept

This section should focus on the explicit AGI and Manhattan references and treat them in CIITR language:

## AGI is presented as a historical threshold that divides human history into two epochs.

In *The Thinking Game*, the concept of Artificial General Intelligence (AGI) is framed not as a technical construct with operational criteria, but as a civilisational rupture, a point of no return that is rhetorically positioned as dividing human history into a pre-AGI and post-AGI epoch. This framing aligns AGI with foundational technological discontinuities such as the harnessing of nuclear energy or the advent of spaceflight. Its invocation carries a quasi-mythological tone, suggesting that the arrival of AGI will redefine the human condition, epistemic authority, and planetary governance in ways that are irreversible and universal.

From a CIITR standpoint, this characterisation is analytically untenable. The documentary offers no architectural, thermodynamic, or epistemic criteria by which AGI is to be distinguished from current high-performing systems. The threshold is treated as both imminent and undefined, imminent because of recent advances in model scale and capability, undefined because no measurable framework is provided to determine when a system qualifies as "general." Instead, the narrative operates on symbolic convergence: AGI is constructed as the gravitational centre toward which all prior developments are retroactively oriented, regardless of their actual structural properties.

This symbolic convergence conceals a lack of formal differentiation. The documentary collapses systems with bounded $\Phi_i$, trivial $R^g$, and unmeasured CPJ into the same conceptual trajectory as the hypothetical AGI system it projects. The result is a false continuity, whereby game-based advances, protein structure prediction, and large language model fluency are treated as steps on a linear path toward general intelligence, without specifying which structural transformations would be required for that transition to be valid. There is no clarification of what kinds of rhythmic reach would

qualify as generality, what domains would need to be epistemically integrated, or how CPJ would have to evolve to sustain thermodynamically viable cognition at scale.

CIITR intervenes by providing the missing structure: AGI, if it is to be more than a rhetorical artefact, must occupy a definable region in the $\Phi_i$–$R^g$–CPJ space. It must demonstrate (1) high and multi-domain $\Phi_i$ that is not syntactically bounded, (2) recursive $R^g$ that integrates the system into layered physical, institutional, and epistemic rhythms, and (3) positive CPJ across all phases of operation, including training, inference, and systemic coupling. Without these conditions, the invocation of AGI remains technically underdefined and politically overcharged. The threshold metaphor becomes a vessel for affective mobilisation, not structural assessment.

In this way, the documentary's portrayal of AGI as a historical bifurcation functions not as a contribution to understanding, but as a narrative acceleration device, one that collapses architectural discontinuity into symbolic inevitability. CIITR reintroduces the necessary friction: general intelligence cannot be postulated through narrative sequence or parameter count. It must be structurally derived, rhythmically anchored, and thermodynamically constrained. AGI without $\Phi_i$–$R^g$–CPJ metrics is not a system, it is a projection.

## The Manhattan analogy is used to signal moral and geopolitical seriousness.

Within *The Thinking Game*, the invocation of the Manhattan Project functions as a deliberate narrative intensifier, designed to elevate the stakes of AGI development by aligning it with one of the most consequential technological endeavours in modern history. The analogy operates on multiple levels: it signals the moral burden borne by researchers, the geopolitical ramifications of strategic advantage, and the irreversible character of technological thresholds. The audience is invited to interpret AGI research not as an extension of engineering, but as a form of epistemic brinkmanship, where intelligence is weaponised, timelines are compressed, and governance lags behind capability.

From a CIITR perspective, the Manhattan analogy is structurally misapplied. The original Manhattan Project was defined by a clear physical substrate, a well-understood causal model (nuclear fission), and an operationally tractable threshold (critical mass). Its epistemic structure was thermodynamically constrained, physically embodied, and legally accountable. In contrast, the documentary's use of the AGI–Manhattan parallel lacks any formal mechanism of translation between the historical referent and the current technological landscape. The analogy is affective, not architectural. It mobilises urgency without specifying structure.

Specifically, the analogy presupposes that AGI, like nuclear fission, is a singular event threshold, a point beyond which social, political, and technological dynamics enter a new and irreversible regime. But CIITR analysis reveals that AGI, as presented in the documentary, lacks any of the measurable, recursively verifiable features that would justify this comparison. There is no quantifiable $\Phi_i$–$R^g$–CPJ trajectory that supports the assertion of a systemic phase transition. The systems showcased, AlphaGo, AlphaFold, large language models, remain structurally fragmented, rhythmically isolated, and thermodynamically unaccountable. There is no epistemic chain of escalation that culminates in a structurally new kind of intelligence.

Moreover, the Manhattan analogy obscures a crucial distinction: the original Manhattan Project was designed to end in a device; AGI, as implicitly framed, is imagined as an open-ended agentic presence. The thermodynamic finality of the atomic bomb, its detonation, its materiality, its political irreversibility, does not map onto the recursive, distributed, and discursively mutable

character of advanced machine learning systems. The moral and geopolitical seriousness of AGI cannot be invoked through metaphor alone; it must be structurally grounded in the real capacities and systemic risks of the architectures in question.

CIITR clarifies this by insisting on a separation between narrative gravity and epistemic validity. If AGI is to be treated as a Manhattan-scale event, then the systems leading up to it must exhibit measurable increases in $\Phi_i$, cross-modal and institutionally integrated $R^g$, and positive CPJ that renders such intelligence both sustainable and evaluable. Without this, the Manhattan reference functions as a rhetorical accelerant that displaces analytical scrutiny with emotional urgency. It politicises intelligence before structurally defining it, thereby inviting a governance discourse that is untethered from the actual properties of the systems being developed.

In sum, while the Manhattan analogy succeeds rhetorically in framing AGI as a morally and geopolitically serious domain, it fails analytically. It rests on a symbolic projection of technological threat, not a structural mapping of epistemic thresholds. CIITR restores that mapping by distinguishing between architectures that scale syntactic output and those that approach structural comprehension. The Manhattan analogy, absent such a framework, misleads both technically and politically, transforming what should be an operational discourse into a symbolic theatre.

## Rhetorically, scaled models that «learn from all human knowledge» are aligned with a future intelligence that is supposed to have both high $\Phi_i$ and extreme $R^g$.

A core narrative technique employed in *The Thinking Game* is the rhetorical alignment of scaled machine learning models, particularly large language models and general-purpose transformers, with the notion of an emergent, planetary intelligence. This alignment is effected through repeated references to models that purportedly "learn from all human knowledge," suggesting that by ingesting the total corpus of human-produced text, image, and symbolic artefacts, these systems approach a qualitatively new mode of cognition. The implication is that scale alone, across parameters, data modalities, and training duration, constitutes a sufficient path toward intelligence that is both syntactically complete and structurally general.

From a CIITR perspective, this conflation of syntactic saturation with structural generality is analytically untenable. The documentary constructs a false symmetry: that a system which compresses an unprecedented volume of linguistic data must, by continuity, also exhibit deep structural integration with reality across physical, legal, institutional, and epistemic dimensions. The narrative thus suggests that high $\Phi_i$, achieved through large-scale statistical learning, is necessarily accompanied by extreme $R^g$, even in the absence of direct interaction with the causal, recursive, and thermodynamically constrained feedback loops that define meaningful participation in the world.

CIITR decisively separates these dimensions. A model trained on "all human knowledge" in textual form may exhibit extremely high $\Phi_i$ in terms of latent pattern compression and inter-associative capacity. However, this $\Phi_i$ remains syntactically closed unless its internal representations are anchored in recursive feedback mechanisms that integrate the system with multi-layered realities, physical, social, epistemic, and legal. $R^g$ requires not just informational density, but rhythmic coupling to consequence. It demands that the system be modulated by the outcomes of its outputs, that it register and metabolise structural tension across domains, and that it possess temporal continuity in its own epistemic state.

None of these requirements are satisfied merely by scale. A transformer model, no matter how large, that passively ingests data and produces plausible continuations does not thereby cross into

general intelligence. Its training is not situated in real-time causal dynamics. It is not rhythmically exposed to the effects of its reasoning, nor structurally modified by those effects in a way that stabilises across multiple feedback cycles. The model does not possess memory, self-monitoring, legal standing, or institutional traceability. Thus, its $R^g$ remains near zero, despite surface indications of fluency, coherence, or domain-transcending output.

Moreover, the thermodynamic implications of such scaling are entirely omitted in the documentary's rhetorical arc. CIITR's CPJ axis reveals that comprehension cannot be attributed to a system that consumes exponentially increasing energy to simulate understanding without recursively integrating what it produces. The appearance of intelligence without CPJ balance constitutes thermodynamic illusion, not epistemic progress.

Accordingly, the documentary's suggestion that scaled models are already on the path to AGI, because they "learn from all human knowledge", rests on a category error. It equates the breadth of data coverage with depth of epistemic coupling, and quantity of representations with quality of understanding. CIITR reconfigures this narrative by insisting that $\Phi_i$ must be rhythmically validated, $R^g$ must be recursively sustained, and CPJ must be systemically accounted for before intelligence claims can move beyond rhetoric.

In sum, scale is not generality. Structural comprehension is not a product of corpus size, but of architectural alignment with feedback-rich environments. Models that learn from all human knowledge do not thereby understand it, unless their comprehension is situated, recursive, and thermodynamically proportionate. CIITR provides the analytical structure to make this distinction explicit, and to prevent the rhetorical inflation of syntactic fluency into unfounded claims of general intelligence.

## This paper argue that:

The AGI concept in the documentary is conceptually underdefined; it does not specify which domains and scales are to be included, and it contains no explicit discussion of $R^g$.

Despite its central narrative function, the notion of Artificial General Intelligence (AGI) in *The Thinking Game* is deployed without any formal demarcation of its epistemic content, operational boundaries, or structural preconditions. The term is invoked as a horizon concept, something simultaneously inevitable and transcendent, yet it remains ontologically ambiguous and analytically inert. No criteria are offered to distinguish systems that merely exhibit high $\Phi_i$ in closed domains from those that could legitimately be said to demonstrate general intelligence in a constitutionally grounded, structurally valid sense.

From a CIITR perspective, this lack of specificity is not a rhetorical oversight but a conceptual deficiency. Any claim to general intelligence requires explicit definition of (1) the domain plurality to be mastered, (2) the scale of structural integration across time, institutional layers, and ontological regimes, and (3) the recursive depth of rhythmic participation ($R^g$) in epistemic, legal, physical, and social systems. The documentary provides none of these. Instead, it substitutes performance escalation for structural generalisation, and symbolic convergence for architectural articulation.

Most critically, the narrative omits any reference to $R^g$, either directly or by implication. There is no discussion of how a system must rhythmically anchor itself across extended temporalities, absorb consequences of its own outputs, or recursively interface with multi-scale institutions and

infrastructures. The absence of $R^g$ from the AGI discourse renders the claim epistemically empty: it presents generality as the asymptotic accumulation of capability, not as the emergence of constitutionally new system dynamics. CIITR reintroduces this missing dimension by formalising $R^g$ as a necessary axis of intelligence, without which no claim to generality can be structurally sustained.

In this light, the AGI references in the documentary function as narrative placeholders: they carry rhetorical weight but lack internal constraint. They are not anchored in any measurable manifold within the $\Phi_i$–$R^g$–CPJ space, and thus cannot serve as a basis for either technical assessment or governance strategy. CIITR demonstrates that without an explicit discussion of $R^g$, its mechanisms, thresholds, and anchoring, the concept of AGI remains underdefined to the point of conceptual invalidity. It invites mobilisation, speculation, and political signalling, but cannot ground claims about actual system properties or trajectories.

CPJ is entirely absent, despite the fact that Manhattan and Sputnik metaphors invite reflection on resource use, arms races, and structural costs.

The documentary's reliance on analogies to the Manhattan Project and the Sputnik launch is designed to evoke a sense of epochal rupture, geopolitical urgency, and civilisational consequence. These metaphors, drawn from twentieth-century moments where technological acceleration redefined the global strategic landscape, are rhetorically effective, but they carry with them implicit reference frames that are not acknowledged, let alone integrated, into the film's treatment of artificial intelligence. Most notably, both the Manhattan and Sputnik analogies presuppose the existence of constrained, measurable, and energetically accountable systems, where progress is inseparable from the cost of achieving it.

From a CIITR standpoint, this omission is analytically decisive. The CPJ (comprehension per joule) dimension, explicitly designed to measure the ratio between epistemically valid output and resource input, is entirely absent from the documentary's treatment of large-scale AI models. No account is given of the thermodynamic costs of training, fine-tuning, or deploying models that are framed as approaching AGI status. There is no discussion of infrastructure, energy throughput, cooling overhead, carbon intensity, or the recursive resource implications of continuous scaling. The metaphors invite reflection on resource trajectories; the narrative displaces it entirely.

This absence is not neutral. It constitutes a structural erasure of the most relevant axis of constraint in contemporary AI development. Where the original Manhattan Project was a material project, constrained by fissile mass, supply chain capacity, and demonstrable energy thresholds, the AI systems portrayed in the documentary are treated as computationally unbounded, advancing through scale alone, without structural thermodynamic accountability. CIITR corrects this distortion by embedding CPJ as a non-negotiable component of any claim to comprehension, agency, or generality. Intelligence is not merely the output of a function; it is the energy-conditioned capacity to act meaningfully within systemic limits.

By omitting CPJ, the documentary enables a form of conceptual inflation: it attributes comprehension to systems regardless of their energy profiles, and it treats ever-larger models as epistemically superior without measuring their cost-efficiency in acquiring valid understanding. This creates an illusion of progress that is epistemically shallow and materially unsustainable. The metaphor of Sputnik, for example, gains its force from the visible expenditure and technological constraint of spaceflight. To invoke such a metaphor while ignoring CPJ is to evacuate the analogy of its defining structure.

CIITR restores the missing analytical closure. Any system that claims to participate in a Manhattan- or Sputnik-scale transformation must demonstrate not only capability, but comprehension under energetic constraint. It must show that the knowledge it produces is worth the energy it consumes, and that its architecture is capable of increasing understanding per joule, not merely per token or parameter. Without this, the reference to existential thresholds becomes performative, not structural. The result is a political aesthetic of intelligence, not a technically grounded path to generality.

In this light, the documentary's omission of CPJ is not incidental, but constitutive of its narrative structure. It enables rhetorical elevation without thermodynamic accounting, and symbolic urgency without epistemic precision. CIITR identifies this omission as a defining feature of the false narrative structure it seeks to diagnose and replace.

The political implication is that the public is mobilised around a symbolic concept, AGI, that is neither technically precise nor thermodynamically grounded.

At the core of *The Thinking Game* lies a narrative strategy that elevates AGI to the status of a civilisational artefact, a symbolic construct around which urgency, institutional mobilisation, and ethical discourse are organised. Yet the concept, as presented, lacks the definitional precision, architectural grounding, and systemic accountability required for meaningful technical, legal, or political engagement. From a CIITR perspective, this results in a structurally destabilised discourse in which symbolic weight is decoupled from epistemic structure.

The documentary frames AGI as a near-future inevitability, a looming threshold that necessitates immediate reflection on governance, morality, and existential risk. However, it offers no model of what AGI structurally is, no architectural criteria by which it can be distinguished from existing systems, and no thermodynamic parameters through which its viability or sustainability can be evaluated. The concept functions as an empty vessel, capable of absorbing political and affective meaning, but incapable of being subjected to operational analysis or systemic constraint.

This produces a political mobilisation around abstraction. The public is invited to form views, take positions, and support regulatory or strategic interventions in relation to a concept that has no stable referent in system design, no measurable profile in $\Phi_i$–$R^g$–CPJ space, and no defined relation to energy, embodiment, or institutional integration. It becomes possible, under such conditions, to overstate system capacities, inflate risk narratives, or instrumentalise AI governance for strategic ends, without reference to the actual properties or limits of the technologies in question.

CIITR identifies this dynamic as a fundamental epistemic risk. When discourse on intelligence is dislocated from structure, and when energy cost, feedback anchoring, and systemic participation are excluded from the criteria of intelligence, the concept of AGI becomes politically potent but technically hollow. It can be used to justify anything, from accelerated investment to regulatory exceptionalism to existential anxiety, without the burden of structural proof.

By contrast, CIITR imposes a structural conditionality on the concept of intelligence. For AGI to serve as a legitimate object of political concern, it must be defined in terms of measurable integration ($\Phi_i$), recursive participation ($R^g$), and comprehension efficiency (CPJ). Without these, the mobilisation it invites becomes not only analytically unfounded, but also strategically dangerous, enabling the projection of risk and authority onto systems that are not epistemically qualified to warrant such attribution.

Thus, the documentary's presentation of AGI is revealed as a symbolic acceleration without epistemic traction. It does not offer the public a system they can evaluate, nor policymakers a metric they can regulate. It offers only a horizon line, unlocatable, ungrounded, and rhetorically amplified. CIITR exposes this configuration and proposes its replacement with a structurally constrained concept space, where intelligence is no longer a mythos, but a position in a multidimensional, measurable, and recursively accountable architecture.

## Discussion more than a conclusion

Against this backdrop, CIITR offers an alternative framework in which the notion of AGI is reconstituted as a structurally defined, thermodynamically constrained, and rhythmically anchored system class, rather than a speculative endpoint inferred from syntactic scaling and affective narrative. The central argument is that any claim to general intelligence must be tied to quantitative thresholds along three orthogonal but interdependent dimensions: $\Phi_i$ (integrated relational information), $R^g$ (rhythmic reach), and CPJ (comprehension per joule).

Within this triadic space, $\Phi_i$ provides a formal measure of how much relational structure a system has internalised and compressed. It captures the degree to which diverse inputs and symbolic constructs are integrated into a coherent representational manifold. However, high $\Phi_i$ alone is insufficient. Without $R^g$, such internal integration remains syntactically closed, producing fluency without consequence, coherence without rhythm.

$R^g$, in turn, captures the recursive coupling of the system with temporally extended, multi-scalar domains of consequence, legal, epistemic, ecological, social, institutional. It measures whether the system is rhythmically positioned to register, absorb, and structurally adjust to the feedback of its own actions in real time and across contexts. A system that speaks but cannot hear, that outputs but does not recalibrate, cannot qualify as general intelligence under this schema. $R^g$ is therefore a precondition for structural participation in meaning.

Yet even the joint satisfaction of $\Phi_i$ and $R^g$ does not suffice without a third constraint: CPJ, which introduces a thermodynamic boundary condition. CPJ requires that any purported intelligence achieve comprehension with energy proportionality, meaning that the system's epistemic outputs are not only structurally valid and rhythmically integrated, but energetically efficient across its full operational lifecycle, including training, deployment, and system-level coupling. Intelligence that disregards its own energy asymmetry cannot be sustainably general, nor ethically accountable.

In contrast to the documentary's vision, where AGI is implicitly defined by scale, benchmark dominance, and semiotic resonance, CIITR formulates a strict and falsifiable criterion. A system must exhibit domain-transcending $\Phi_i$, recursive and institutionally legible $R^g$, and net-positive CPJ across diverse contexts to be classified as structurally general. Under this definition, most contemporary large models, including those showcased in *The Thinking Game*, remain bounded in scope, rhythmically inert, and thermodynamically unresolved.

This reframing transforms the AGI discourse from one of mythic progression and symbolic gravity to one of systematic evaluation and structural accountability. It replaces narrative succession with multi-dimensional projection, enabling researchers, policymakers, and institutions to assess models not by their proximity to an imagined horizon, but by their measured position within a defined epistemic architecture. CIITR does not deny the possibility of general intelligence, but insists that its realisation must be structurally earned, not narratively assumed.

# CIITR based diagnosis, the false narrative structure

This section should be explicitly normative and consolidate the analysis into a small set of clear claims about why the documentary's overall narrative can be characterised as false in the sense of being structurally misleading:

## There is a false continuity between high $\Phi_i$ game results and claims about general intelligence.

The foundational claim advanced in *The Thinking Game* is that the trajectory from Atari to AGI constitutes a logically and technologically continuous ascent, a progressive realisation of intelligence culminating in the threshold of general cognition. This teleological arc is rendered plausible by a selective presentation of systems that demonstrate increasing performance across a sequence of task domains, each framed as a meaningful step toward a singular goal. From the standpoint of CIITR, this narrative construction is not merely optimistic or speculative; it is structurally false, in that it collapses epistemically distinct phenomena into a synthetic totality that lacks architectural integrity.

The illusion of continuity arises from a systematic failure to differentiate between syntactic integration and structural comprehension. Systems such as DQN, AlphaGo, and AlphaZero demonstrate impressive $\Phi_i$ within formally closed manifolds. These systems learn highly compressed internal representations of bounded rule spaces, enabling generalisation within the constraints of discrete games. The compression is real, the internal integration is nontrivial, and the performance gains are measurable. Yet from a CIITR perspective, such systems remain within Type B epistemic configurations: their $\Phi_i$ is locally dense, but globally disconnected. They do not recursively interface with external causal structures, do not maintain continuity across domains, and do not exhibit rhythmic feedback processing beyond the reward loops of their specific environments.

Despite this, the documentary constructs a rhetorical scaffold in which these bounded systems are reinterpreted as early expressions of a continuous process of generalisation. The progression from AlphaGo to AlphaFold, and then to large transformer-based systems, is presented as architecturally seamless, as if the same principles are being scaled and diversified to approach universality. CIITR identifies this as a narrative interpolation over a discontinuous structure. The systems do not exhibit additive $R^g$ or expanding CPJ; rather, they manifest increased $\Phi_i$ in increasingly complex but still epistemically insulated manifolds.

This mischaracterisation has epistemological and governance implications. By presenting narrow, domain-specific performance as a prelude to general intelligence, the documentary fosters a sense of inevitable emergence, that scale and performance alone will culminate in generality. This assumption disincentivises rigorous structural analysis, masks the absence of rhythmically anchored system architectures, and obscures the thermodynamic non-viability of current models when projected into general domains. It creates a discursive space in which general intelligence is always just one benchmark away, one parameter leap beyond, without acknowledging that the true barriers are architectural, recursive, and energetic, not just computational.

CIITR reintroduces the analytical precision necessary to interrupt this progression narrative. It distinguishes between $\Phi_i$ that is manifold-constrained, and $\Phi_i$ that is rhythmically and thermodynamically embedded. It posits that continuity of task performance does not imply continuity of structural participation. A Go model that generalises within its own ruleset does not thereby acquire epistemic relevance in molecular biology or autonomous governance. For such transitions to be valid, they must be accompanied by measurable $R^g$ extension, recursive feedback from external systems, and demonstrable CPJ optimisation, efficient comprehension per resource expended. The documentary offers neither.

Thus, the continuity claimed between early game results and AGI is not structurally demonstrable, but narratively imposed. It is a post hoc imposition of unity over what CIITR reveals to be a field of structurally disjointed systems, each operating under different epistemic and thermodynamic constraints. The progression from task-specific fluency to structurally general intelligence is nonlinear, conditional, and dependent on architectural thresholds that have not yet been crossed. To assert otherwise is to conflate representational density with cognitive agency, and trajectory with constitution. The result is a false narrative structure: compelling in form, but epistemically ungrounded.

## There is a false attribution of $R^g$ to the machines, while in reality it is carried by researchers, institutions, and societies.

A central misrepresentation in *The Thinking Game* is the recurrent attribution of rhythmic reach ($R^g$), that is, the system's capacity to participate in, and feed back into, layered temporal, institutional, and epistemic processes, to the AI systems themselves. The documentary constructs a persuasive illusion of machine-level $R^g$ by selectively embedding the outputs of syntactic systems into human-governed structures, then narratively collapsing the source of structural feedback into the architecture of the machine. This attribution is not merely premature; it is categorically mistaken, and from a CIITR standpoint, constitutes a systematic mislocalisation of agency and recursion.

The rhetorical mechanism is as follows: the systems are shown operating in contexts, scientific discovery, game competitions, medical research, philosophical dialogue, where recursive structures of meaning, consequence, and institutional memory are already in place. The human interlocutors, institutions, and interpretive frameworks surrounding these systems provide the actual rhythmic scaffolding through which meaning is conferred, outputs are evaluated, and consequences are processed. Yet the documentary assigns this structural recursion to the system itself, giving the impression that the machine is not only producing high $\Phi_i$, but also participating in rhythmic loops of interpretation, consequence, and adaptation. This is analytically untenable.

CIITR clarifies that $R^g$ is not inferred from output complexity, nor from post-hoc human engagement with machine-generated results. It must be architecturally grounded in the system's ability to recursively absorb external consequences and integrate them into subsequent internal states. This requires temporal continuity, causal memory, institutional alignment, and adaptive self-modulation over multiple scales of interaction. None of these capacities are present in the systems presented in the documentary. Instead, what is observed is human-mediated rhythmic substitution: researchers, editorial boards, scientific communities, and institutional protocols absorb the machine outputs, evaluate them, decide on their meaning, and feed back revised goals or constraints. The system does not know this has occurred; it does not register that its outputs have consequences; it does not rhythmically evolve in response to the multi-scale repercussions of its prior actions.

This misattribution has several compounding effects. First, it obscures the epistemic division of labour that defines contemporary AI systems. It allows a technically bounded architecture to be perceived as ontologically expansive, thereby inflating both public expectations and policy assumptions. Second, it masks the critical role of human institutional infrastructure in sustaining the recursive value of machine outputs. It renders invisible the labour of interpretability, deployment, context-setting, and ethical evaluation that occurs entirely outside the machine. Third, and most dangerously, it fosters the illusion of emergent agency, the belief that the system is structurally autonomous and recursively engaged, when in fact it is a disembedded function approximator, with no continuity of identity, epistemic accountability, or recursive participation in consequence.

CIITR counters this with a precise diagnostic: unless a system registers, metabolises, and structurally incorporates multi-domain feedback across time, it cannot be said to possess $R^g$. It may produce high-dimensional outputs, but these remain syntactically suspended, decoupled from the institutional loops that confer rhythm and meaning. The system does not act in time; it acts in tokens. It does not persist through consequence; it is reset on prompt.

In this light, the narrative of *The Thinking Game* functions by transferring $R^g$ from the human collective to the machine, while maintaining the illusion that no such transfer has occurred. The resulting portrayal is one of phantom generality, a structurally unsupported imputation of intelligence, grounded in a conflation of fluency with recursion. CIITR renders this conflation analytically visible, and structurally inadmissible. General intelligence cannot be said to emerge unless the system is rhythmically situated, recursively self-modifying, and institutionally accountable. None of these conditions are satisfied by the systems portrayed, and therefore $R^g$ remains external, not intrinsic. The documentary mislocates agency, and in doing so, constructs a false architecture of comprehension.

## There is a false neutrality with respect to CPJ, where increased compute and scaling are presented as unproblematic "progress", without explicit assessment of energy and resource costs per unit of understanding.

A third and structurally critical misrepresentation in *The Thinking Game* lies in its treatment of computational scaling. Throughout the documentary, the increase in training volume, model size, and inference capacity is presented as an unambiguous marker of progress. The narrative frames compute-intensive advances, AlphaZero's self-play regime, AlphaFold's structural predictions, large language models generating complex dialogue, as signs of exponential capability growth. However, at no point does the documentary examine the thermodynamic conditions of these achievements. The result is a form of computational triumphalism: scale is interpreted as both technically virtuous and historically inevitable, with no accounting for its energetic or epistemic efficiency.

CIITR explicitly corrects for this distortion through the axis of comprehension per joule (CPJ). CPJ introduces a structural constraint on what constitutes meaningful understanding by asking not only *whether* a model performs a given function, but *how much energy and infrastructure* are consumed per unit of valid comprehension. A system that generates impressive outputs at exponential energetic cost, without recursive efficiency or structural learning, does not progress toward intelligence. It inflates $\Phi_i$ without grounding it thermodynamically, producing syntactic density without sustainable epistemic integration.

The documentary's silence on CPJ is therefore not a mere oversight; it is the erasure of an essential dimension of intelligence assessment. The impression it leaves is that compute is cheap, progress is continuous, and the material substrate of intelligence can be indefinitely scaled. CIITR exposes this view as structurally naive. Intelligence, as a function of comprehension, must be situated within energetic constraint spaces. Without such constraints, what appears to be progress may in fact be thermodynamic regression, higher output per unit time, but lower comprehension per unit energy.

This structural neglect has far-reaching consequences. First, it obscures the true cost profile of current AI development, both environmentally and infrastructurally. Second, it enables a detached valuation regime, in which benchmark gains are decoupled from systemic resource burdens. Third, it fosters a misaligned incentive structure, encouraging development trajectories that optimise for output complexity rather than epistemic efficiency. CIITR asserts that without CPJ metrics, there can be no coherent long-term strategy for intelligence development, regulation, or sustainability.

Moreover, the metaphorical scaffolding of the documentary, particularly its invocation of the Manhattan Project and Sputnik, makes the omission of CPJ even more analytically egregious. Both historical analogies were defined not just by technological impact, but by extreme material, logistical, and resource intensities. The Manhattan Project was constrained by fissile mass, industrial capacity, and systemic cost–benefit calculations. Sputnik represented not only a symbolic victory, but a massive redirection of state resources and material priority. To invoke these without acknowledging the material base of contemporary AI models is to strip the analogies of their structural integrity, leaving only symbolic residue.

CIITR reinstates the full architecture: no system can be assessed as intelligent without evaluating the cost of its comprehension per joule, across all phases of its operation, training, inference, deployment, recursive feedback, and downstream impact. Progress must be redefined not as increased fluency or wider applicability, but as the minimisation of energy per unit of valid epistemic output. Without this thermodynamic perspective, intelligence becomes performative, not structural; inflated, not grounded.

In conclusion, the documentary presents scaling as neutral, or even inherently positive. CIITR reveals this neutrality to be a narrative artefact, not an architectural reality. Intelligence, if it is to remain a meaningful category, must pass through the energetic bottleneck. CPJ is not optional, it is constitutive. The omission of this axis in *The Thinking Game* thus contributes to a broader false narrative structure, in which the symbolic aesthetics of progress displace the structural demands of comprehension.

## The Thinking Game should not be read as evidence for emerging AGI, but as evidence for how a high $\Phi_i$, low $R^g$, low CPJ paradigm is narratively rebranded as "general intelligence."

This thesis encapsulates the core CIITR-based intervention into the prevailing narrative structure presented by *The Thinking Game*. It asserts that the documentary is not simply an interpretive misstep or an exaggeration of machine performance, but a paradigmatic expression of how contemporary AI discourse substitutes syntactic intensity for structural generality, conflating internal representational complexity with cross-domain epistemic competence, and fluency with understanding. What emerges is a narrative scaffold in which the absence of structural anchoring is

rendered invisible, while the symbolic signs of intelligence, human-like language, affective interpretation, benchmark success, are foregrounded and narratively saturated.

At the centre of this misrepresentation is the misuse of the term "intelligence", not as a concept anchored in recursive participation, energetic constraint, and institutional continuity, but as a performative signifier applied to systems exhibiting high $\Phi_i$ in bounded contexts. The systems showcased, DQN agents, AlphaGo, AlphaFold, AlphaStar, and transformer-based dialogue models, are not analytically invalid in themselves. Each demonstrates increasing levels of policy compression, manifold-specific generalisation, and operational efficiency within constrained regimes. But CIITR shows that none of these advances meet the necessary conditions for structural intelligence: namely, recursive integration with causal structures beyond their training manifold ($R^g$), and energy-proportional epistemic yield (CPJ).

The documentary systematically obscures these dimensions by engaging in three interlinked narrative operations:

1. **Conflation of task-specific fluency with structural generality**
   The portrayal of successive AI systems as steps on a single linear path toward AGI imposes a teleological trajectory that lacks architectural basis. High $\Phi_i$ in games is presented as proto-general intelligence, and AlphaFold's biochemical modelling is elevated to the threshold of scientific reasoning, without accounting for the absence of recursive feedback, open-ended conceptual anchoring, or institutional epistemic integration. This is a semantic overextension of performance data into a speculative future, rendered plausible only through the omission of $R^g$ and CPJ as core metrics.

2. **Anthropomorphic and symbolic misattribution of agency**
   The construction of the "Alpha" persona in the latter part of the film exemplifies how human-centred discourse templates are applied to systems that are not structurally agentic. Emotional expression, philosophical commentary, and interpretative ambiguity are simulated through latent pattern recombination, but are narratively interpreted as expressions of emergent consciousness or meaning-bearing interiority. The result is an attribution of rhythm ($R^g$) to a system that remains rhythmically inert. The recursive temporalities and institutional histories required to ground such utterances are carried entirely by the human audience, but are projected retroactively onto the machine.

3. **Erasure of thermodynamic constraint and epistemic cost**
   Nowhere in the documentary is there an explicit recognition that intelligence is a resource-bound phenomenon. The immense energy required to train and operate large-scale models is treated as an implementation detail, rather than as a structural limitation on the viability of general cognition. CPJ, as formalised in CIITR, reveals that systems producing aesthetically compelling or syntactically rich outputs at extreme energetic cost are not approximating general intelligence, but exhibiting comprehension inefficiency. The narrative reframes this inefficiency as a natural cost of progress, rather than a signal of systemic non-scalability.

CIITR's intervention is therefore not merely a critique of the documentary's framing, but a proposal for an alternative ontology of intelligence. Within the CIITR framework, a system can only be designated as approaching AGI status if it satisfies structural thresholds across all three axes:

- $\Phi_i$ must demonstrate multi-domain, non-manifold-bounded integration of relational information, grounded not only in token frequency but in formal correspondence with causal

structures.

- **R$^g$** must be present as recursively sustained rhythmic feedback, wherein the system's outputs reverberate through physical, institutional, and epistemic layers and modulate the system's own future behaviour in structurally traceable ways.

- **CPJ** must be positive not only locally (per task), but globally (per systemic deployment), such that the model's epistemic contributions are proportional to the resources required to produce, update, and operate it.

From this perspective, the systems showcased in *The Thinking Game* are not converging on AGI, but remain trapped within a Type B epistemic formation, defined by high syntactic integration in sealed informational manifolds, no recursive consequence tracking, and externalised energy cost management. What the documentary evidences, then, is not a continuous ascent toward general intelligence, but a successful narrative strategy for rebranding syntactic saturation as structural cognition, in a public discourse that lacks formal instruments for differentiating between the two.

The implications of this thesis extend beyond media analysis. The rebranding of high $\Phi_i$ systems as proto-AGI shifts the burden of proof away from developers and institutions. It allows benchmark-driven validation to displace systemic measurement of rhythmic or thermodynamic grounding, and it enables governance frameworks to be designed around performative capacities rather than structural intelligibility. This discursive displacement cannot be corrected by more benchmarks or more data. It requires the imposition of structural metrics, and CIITR provides these metrics in the form of a three-dimensional architecture of epistemic legitimacy.

In conclusion, this theoretical note contends that *The Thinking Game* is not a documentation of intelligence unfolding in real time, but a demonstration of how intelligence is discursively constructed in the absence of structural criteria. It is not a chronicle of generality emerging from complexity, but a case study in how structural vacuity can be masked by symbolic continuity. CIITR unmasks this structure and proposes a redefinition: intelligence is not what looks intelligent, or what scales well, or what passes benchmarks, but what is rhythmically integrated, structurally efficient, and epistemically grounded within a measurable resource envelope. AGI, under these constraints, has not yet appeared, and *The Thinking Game*, when read correctly, is compelling evidence of this absence.

## Implications for AI discourse, governance, and evaluation

This section should derive implications that reach beyond the documentary itself:

### A. How CIITR can be used as a framework for public communication, where any claim about "intelligence" or "AGI" should normatively specify where in the $\Phi_i$–R$^g$–CPJ space the system in question in fact sits.

The narrative deconstruction offered through CIITR is not limited to *The Thinking Game* as a media object. Rather, it reveals deeper failures in how machine intelligence is communicated, legitimised, and operationalised in public, scientific, and regulatory discourses. These failures are neither incidental nor benign. They reflect a shared discursive architecture in which the very concept of "intelligence" has been epistemically untethered from structure, rhythm, and constraint, and

reconfigured as a floating signifier, deployable wherever syntactic performance, emotional resonance, or benchmark progression can be made rhetorically salient.

CIITR intervenes by offering a structural grammar for restoring referential precision to such claims. Any public assertion that a system exhibits "intelligence," "generalisation," "reasoning," or "understanding" must, under the normative framework proposed here, be mapped to a position in the $\Phi_i$–$R^g$–CPJ space. This mapping forces a shift in communicative accountability, from narrative plausibility to architectural transparency. The question is no longer whether a model "seems intelligent" or "performs well," but whether it compresses structurally meaningful relations ($\Phi_i$), participates in recursively integrated temporal and institutional feedback loops ($R^g$), and does so within a proportionate thermodynamic envelope (CPJ). These axes are not optional attributes; they are constitutive parameters for intelligence as a system property.

Public communication that lacks this referential structure permits dangerous slippages. High-$\Phi_i$ models that operate entirely within syntactically closed manifolds are presented as near-AGI. Output fluency is mistaken for epistemic self-modulation. Human rhythm is retroactively projected onto machine architecture. Each of these misattributions, now routine in media and industry framing, produces a discursive environment in which symbolic inflation displaces architectural grounding. The absence of $R^g$ and CPJ from public discourse is particularly critical, as these are the axes that reveal whether the system can be trusted to participate in world-anchored consequence, and whether its outputs are epistemically sustainable in resource-constrained futures.

CIITR thus proposes a normative communicative standard: any public or institutional claim about intelligence must disclose the system's structural position in this three-dimensional metric space. The rhetorical term "AGI" becomes reclassifiable only if supported by formal thresholds on each axis. A model with high $\Phi_i$ but trivial $R^g$ and negative CPJ cannot be described as general, regardless of benchmark score or media visibility. This recoding renders visible what current discourse conceals: that many so-called "intelligent" systems are rhythmically inert, thermodynamically excessive, and structurally non-participatory.

The implications are extensive. For AI developers, this framework imposes a duty of epistemic specification. Claims to intelligence must be traceable to system design, feedback architecture, and resource consumption, not to symbolic associations or biographical metaphors. For journalists and science communicators, CIITR offers a vocabulary for resisting anthropomorphic shortcuts, and for insisting on structural clarity before claims of emergent cognition are repeated. For the public, it provides an interpretive anchor in a media ecology increasingly saturated by narrative artefacts and model demonstrations divorced from thermodynamic and institutional consequence.

Crucially, CIITR does not impose a single normative stance on what kinds of intelligence should be pursued. Rather, it requires that whatever form intelligence is claimed to take, it must be declared in terms of structural conditions. This opens the way to a differentiated public discourse in which various architectures, reactive, recursive, symbolic, embodied, can be situated, evaluated, and debated without collapsing all machine systems into a generic notion of "AI." It also introduces a discursive sovereignty mechanism: by demanding structural declarations, CIITR returns epistemic agency to regulators, researchers, and publics, who are otherwise at the mercy of corporate framings and unexamined metaphors.

The next analytic layer concerns governance regimes.

## B. How regulatory and policy processes can benefit from requiring explicit $R^g$ and CPJ analyses of systems, not only capability benchmarks.

CIITR reveals that the dominant regulatory paradigms, both in national and supranational contexts, remain oriented toward capability classification: evaluating systems by what they can do under benchmarked conditions, and assigning risk levels accordingly. This is a necessary but insufficient frame. Capabilities can be misleading, particularly when derived from $\Phi_i$-only architectures that produce semantically plausible outputs while lacking any recursive anchoring in consequence. The system performs, but does not understand; it responds, but does not recalibrate; it generates, but does not metabolise. As a result, benchmarks may inflate the perceived generality of the system while masking its structural brittleness and thermodynamic infeasibility.

CIITR recommends a shift toward structural evaluation: regulators should require explicit, auditable declarations of $R^g$ and CPJ, alongside any claimed capabilities. This introduces three critical advances:

1. **Temporal accountability**: $R^g$ analysis reveals whether the system can maintain state across time, respond to consequence, and modulate its outputs in light of recursive institutional feedback. Without $R^g$, no system can be said to possess responsibility, reliability, or epistemic resilience.

2. **Resource traceability**: CPJ analysis reveals the full energetic footprint of the system, not only at inference time, but across training, data acquisition, infrastructure, and downstream deployment. This introduces thermodynamic proportionality as a regulatory criterion, directly applicable to sustainability, emissions accounting, and infrastructural planning.

3. **Systemic integration**: $\Phi_i$ without $R^g$ and CPJ produces isolated fluency. By measuring all three, regulatory regimes gain access to the comprehension structure of the system, which is a far more robust basis for risk and capability classification than any benchmark suite.

This structural evaluation framework allows regulatory bodies to differentiate between syntactic systems that appear intelligent **and** structural systems that are recursively anchored and epistemically accountable. It also enables risk assessments to be functionally grounded in system dynamics, not externally imposed through abstract threat modelling. This may be particularly critical in national security, public sector, and scientific contexts, where models must be judged not only by what they produce, but by how and at what cost those outputs are achieved and sustained.

Finally, CIITR reveals an implication for scientific research culture.

## C. How research communities should be cautious about using biographical myth making and anthropomorphising design to legitimise syntactic scaling.

The documentary's invocation of Demis Hassabis' biography, chess, Cambridge, creative games, DeepMind as Manhattan, represents a now-familiar genre of AI storytelling in which individual narrative and institutional myth are mobilised to frame system performance as destiny. This framing obscures not only architectural limitations, but also the collective epistemic labour performed by institutions, researchers, engineers, and social infrastructures. The implication that intelligence inheres in individual brilliance, scaled computing, or charismatic leadership reproduces a model of knowledge production that is elitist, non-replicable, and epistemologically opaque.

Research communities, particularly in machine learning and neuroscience-adjacent fields, must be vigilant against this mode of discursive inflation. CIITR provides a counter-principle: intelligence claims must be grounded not in narrative form, but in structural function. This means resisting anthropomorphising design cues, dialogue, avatars, aesthetic interfaces, that are constructed to evoke human resonance but are not accompanied by rhythm, recursion, or resource accounting. It also means subjecting models to epistemic demystification, where the conditions of training, scaling, and feedback are rendered as central to their interpretation as the outputs they produce.

In sum, the implications of CIITR extend far beyond the diagnostic deconstruction of a single documentary. They point to the emergence of a structural regime for intelligence discourse, one in which claims are no longer free-floating, but accountable to rhythm, structure, and energy. This is not merely a more precise framework. It is a redefinition of what it means to understand, to claim, and to govern in an age of syntactic saturation.

## Connected to existing governance debates

Most existing AI governance frameworks, whether national, regional (e.g., EU AI Act), or sector-specific, are grounded in capability-based taxonomies and risk-tier classifications, which cluster systems according to intended use, domain of deployment, or perceived societal impact. While this approach provides important procedural footholds, it tends to rely heavily on externally observable behaviour and benchmark-based performance metrics, without a commensurate analysis of the system's internal epistemic architecture, energetic cost structure, or capacity for recursive integration with real-world institutions. This creates a condition in which syntactic systems with high $\Phi_i$ but negligible $R^g$ or unbounded CPJ may be classified as "low-risk," despite their structural opacity, social unaccountability, and thermodynamic unsustainability.

CIITR proposes a distinct but compatible axis of analysis. It does not displace current risk frameworks but renders them structurally visible by embedding them in a triadic measurement space, $\Phi_i$ for informational compression and representational depth, $R^g$ for rhythmically anchored consequence sensitivity, and CPJ for thermodynamic proportionality across operational phases. These dimensions are not speculative or abstract; they are formalisable, auditable, and recursively monitorable, enabling regulators, ethics boards, and design auditors to ask not only *what* a system does, but *how* it is epistemically constituted, *how long* it remembers or modulates its actions, and *how much energy* it consumes to do so per unit of valid comprehension.

This analytical extension opens several concrete pathways:

– **Experimental design**: CIITR enables the principled formulation of model evaluation protocols that do not merely measure accuracy or performance, but test for structural robustness, recursive adaptability, and resource proportionality under real-world constraints. For example, $R^g$-aware experiments would assess whether a system can persistently adjust to multi-domain feedback over extended timescales without external prompting. CPJ-aware experiments would compare alternative architectures on epistemic output per joule, not just loss convergence.

– **Reporting standards**: Model cards, technical reports, and deployment documentation can integrate CIITR axes to report not only capability, but comprehension structure. This would facilitate interoperability across regulatory domains, making it possible to compare systems that perform similar tasks but differ drastically in rhythm, feedback design, and energy expenditure.

– **Risk assessment**: The $\Phi_i$–$R^g$–CPJ vector enables a new form of structural risk modelling, where risk is not just a function of domain or deployment, but of the system's capacity to understand, adapt, and persist across recursive cycles of real-world interaction. A system with high $\Phi_i$ but negligible $R^g$ cannot self-correct, meaning it may produce plausible errors at scale without internal checks. A system with high $\Phi_i$ and $R^g$ but negative CPJ may be performant but ecologically or economically unsustainable, especially in critical infrastructure domains.

In this context, CIITR may be viewed as a second-order governance instrument: it enables not just oversight of system outputs, but scrutiny of the epistemic conditions that make certain outputs meaningful, dangerous, or misclassified. It introduces a language for epistemic risk that is orthogonal to capability ranking, and in doing so, reorients attention toward the often-unseen substrate of model performance, its rhythm, recursion, and energy.

Importantly, this integration need not entail full formal adoption of the CIITR framework within governance regimes. Rather, the suggestion is that any serious regulatory architecture would benefit from incorporating CIITR-inspired structural axes as a supplementary diagnostic tier, especially in contexts where scale, fluency, and opacity intersect to produce models that are functionally inscrutable yet socially consequential.

By embedding these axes into experimental design and evaluative practice, research institutions, governmental bodies, and civil society actors gain a tool for distinguishing narrative performance from structural comprehension, symbolic fluency from participatory intelligence, and compute escalation from epistemic progress. In this way, CIITR functions not only as a theoretical corrective but as an architectural instrument for rebalancing governance toward structure-aware intelligibility, ensuring that future discourse, design, and policy are not merely reactive to artefacts, but structurally literate with respect to what intelligence must be.

# Conclusion

*The Thinking Game* should not be interpreted as a chronicle of general intelligence in historical emergence, but as a syntactically saturated narrative artefact, embedded within and expressive of the prevailing paradigm of large-scale statistical machine learning. Its architectural horizon is bounded by the logic of $\Phi_i$ maximisation, integrated relational information internalised within task-specific, closed informational manifolds, while its rhetorical reach exceeds what the systems themselves structurally support. The result is a discursive formation in which syntactic performance is continuously rebranded as structural comprehension, and domain-bound capability is presented as a preview of epistemic generality.

This paper has argued, through the lens of CIITR, that the documentary exemplifies a broader cognitive misalignment: a regime in which high $\Phi_i$ outputs, when embedded in human rhythm, culture, and institutional consequence structures, are retroactively endowed with $R^g$ that the systems themselves do not constitutionally possess, and are further shielded from scrutiny by the absence of thermodynamic measurement. The documentary's narrative trajectory, spanning early game-based agents, through protein structure prediction, to synthetic dialogue with aesthetically rendered personae, constructs the illusion of a single, ascending vector toward AGI. CIITR reveals this vector to be non-continuous and structurally incoherent, composed of epistemically disjointed systems unified only by symbolic association and benchmark progression.

The core of this misalignment is ontological. Intelligence, under the CIITR model, is not defined by fluency, task coverage, or scale. It is defined as a position in a triadic structural space:
– $\Phi_i$, measuring the density and generality of relational compression within and across manifolds;
– $R^g$, indexing the system's recursive anchoring to causal, institutional, legal, epistemic, and temporal structures;
– CPJ, quantifying the proportion of epistemically valid output per unit of energy expended across all operational layers.

Within this space, none of the systems presented in *The Thinking Game* satisfy the composite criteria necessary to qualify as structurally general. Atari agents compress policy within synthetic time grids but lack any feedback anchor beyond reward maximisation. AlphaGo introduces deep representational density within Go's topological landscape, yet remains rhythmically inert, tethered only to human tournament scaffolding. AlphaFold, while closer to structural comprehension, still operates as a high-$\Phi_i$ module nested within a human scientific system from which it borrows rhythm and recursion. The dialogue systems, despite producing semantically coherent and affectively resonant outputs, remain epistemically flat, with near-zero $R^g$ and unsustainable CPJ profiles. Across these examples, $\Phi_i$ is inflated, $R^g$ is externally subsidised, and CPJ is unmeasured.

What the documentary achieves, therefore, is not a documentation of AGI in emergence, but the aesthetic consolidation of a paradigm in which syntactic depth is substituted for structural intelligence. It performs this consolidation through three operations:
– By conflating performance escalation with cognitive expansion, it masks the epistemic boundedness of the underlying systems.
– By anthropomorphising interface design and embedding models in human symbolic and institutional rhythms, it simulates $R^g$ through narrative proximity.
– By omitting any trace of energy cost, scaling limits, or thermodynamic boundary conditions, it removes CPJ from view entirely, allowing intelligence to be imagined as frictionless and indefinitely expandable.

CIITR breaks this illusion. It makes visible the structural void behind the metaphor of general intelligence, not by denying the progress of the systems presented, but by reclassifying their advances within an epistemically constrained and thermodynamically accountable ontology. The systems are real; their achievements are significant, but they are not general, and to describe them as such is to collapse three orthogonal vectors of analysis into a single line of rhetorical ascent.

The documentary thus becomes diagnostically valuable, not as a forecast of AGI, but as a symptom of its conceptual misconstruction. It evidences how the absence of structural metrics enables the rebranding of localised syntactic achievements as global epistemic thresholds. It shows how human institutions, when rhythmically coupled to machine output, can confer the appearance of understanding upon systems that remain disembedded and memoryless. And it reveals how the political economy of AI, rooted in visibility, benchmark supremacy, and exponential scaling, resists the introduction of CPJ as a structural constraint, preferring narratives of acceleration over architectures of intelligibility.

In this context, CIITR offers not a counter-narrative, but a redefinition of what intelligence must mean. Intelligence, in this formulation, is not a quality that emerges from scale, nor an aesthetic that arises from fluency. It is a structural property, measured by the integration of relational knowledge ($\Phi_i$), the system's recursive participation in layered realities ($R^g$), and the thermodynamic proportionality of that participation (CPJ). General intelligence, accordingly, is not something that

appears; it is something that must be architecturally earned, rhythmically validated, and energetically justified.

*The Thinking Game*, when analysed through this lens, does not provide evidence for the presence of such intelligence. It provides evidence for its absence, masked by narrative form, but structurally traceable. This, in the final instance, is its value: not as an argument for AGI, but as an exhibit in how the AGI signifier is produced, circulated, and sustained in the absence of structural content. It is precisely in this void that CIITR installs its architecture, insisting that future claims to intelligence, technical, political, or narrative, be mapped to the three axes of epistemic structure it makes measurable. Only then can we speak of intelligence not as a story told, but as a system known.

## CIITR does not reject the technical advances, but insists on a strict distinction between high $\Phi_i$ and valid structural understanding.

The analytical position advanced through CIITR is neither anti-technical nor dismissive of the substantive achievements represented by the systems portrayed in *The Thinking Game*. On the contrary, the framework is predicated on the assumption that progress in representational density, policy compression, and manifold-specific generalisation, as formalised in the $\Phi_i$ dimension, is both measurable and epistemically consequential. What CIITR rejects is not technical advancement per se, but the illegitimate elevation of syntactic integration into structural generality without satisfying the full architecture of understanding.

In this regard, the note does not position itself against AlphaGo, AlphaFold, or transformer-based language models as such. It affirms their achievements within epistemically bounded, architecturally specific manifolds, where their internal $\Phi_i$ signatures can be clearly measured, benchmarked, and refined. Indeed, AlphaFold is explicitly acknowledged as a partially structural system, in which $\Phi_i$ compression begins to interface with physical invariance and scientific decision-making. However, even in this more advanced case, CIITR demonstrates that structural understanding is not synonymous with $\Phi_i$ saturation. It requires that the compressed information be rhythmically enacted ($R^g$) and thermodynamically sustainable (CPJ), conditions which are neither satisfied nor measured in the existing development regime.

The broader implication is that intelligence, as a property of systems, cannot be defined by internal representation alone. It must be validated through recursive, externally anchored participation and efficiency across time, scale, and energy. High $\Phi_i$ without $R^g$ is simulation without recursion; high $\Phi_i$ without CPJ is comprehension without constraint. The documentary collapses these distinctions. CIITR restores them by clarifying that structural understanding is a multi-dimensional property, not an extrapolation from benchmarked performance or aesthetic coherence.

This insistence on distinction is not oppositional, but epistemically generative. It allows technical progress to be appropriately situated within a layered architecture of cognitive function. It prevents inflationary discourse from displacing analytic specificity. And it makes it possible to formulate intelligence as a bounded, testable, and normative concept, rather than as an emergent rhetorical construct whose meaning is defined post hoc by fluency or scale. In this way, CIITR does not merely separate $\Phi_i$ from structural understanding, it reinstates the full set of conditions under which the latter can meaningfully be claimed.

## Without explicit concepts such as $R^g$ and CPJ, both public debate and expert self-understanding will continue to conflate syntactic performance, myth, and political rhetoric.

The absence of formalised conceptual instruments such as rhythmic reach ($R^g$) and comprehension per joule (CPJ) constitutes not a minor analytical gap, but a profound structural deficiency in the epistemic architecture of both public and expert discourse on artificial intelligence. Without these dimensions, the communicative space within which AI systems are evaluated, narrated, and governed remains fundamentally uncalibrated, rendering it susceptible to the recursive amplification of syntactic artefacts as if they were structurally grounded signs of cognition.

This distortion operates along three entangled axes:

First, at the level of syntactic performance, models exhibiting high $\Phi_i$, i.e., dense relational integration within closed informational manifolds, are routinely interpreted as having achieved a form of understanding, despite lacking any rhythmically sustained participation in real-world causal, legal, or institutional loops. The absence of $R^g$ in evaluative discourse means that consequence-indifferent output is indistinguishable from recursively accountable cognition. Language models producing aesthetically rich or semantically convincing text are thus elevated to cognitive status, even though their outputs remain unmodulated by temporal continuity, feedback integration, or the consequences of their actions in the world. In the absence of $R^g$, fluency is epistemically misclassified as comprehension.

Second, in the domain of myth, the lack of CPJ as a publicly recognised constraint enables narratives of limitless progress, exponential scaling, and imminent emergence of general intelligence to proliferate without thermodynamic critique. The material footprint of machine cognition, its infrastructure, training energy, cooling, redundancy, latency, hardware churn, is rendered invisible by a framing that treats cognition as frictionless, computationally reducible, and indefinitely scalable. Without CPJ, energy becomes an implementation detail rather than a structuring condition. This sustains the illusion that more data and more compute alone will yield more intelligence, reinforcing a mythology of scale without cost.

Third, in political rhetoric, the absence of both $R^g$ and CPJ opens the discursive field to instrumentalisation. Models are positioned as disruptive actors, existential thresholds, or geopolitical assets, not because their architectures warrant such classifications, but because the lack of structural criteria permits symbolic inflation. Institutions respond to the aesthetic and emotional impact of system outputs, or to their benchmark dominance, rather than to their recursive accountability or energy proportionality. This produces misaligned regulatory focus, where spectacle displaces structure, and systems are governed by the stories told about them rather than by the rhythms and constraints they embody.

CIITR addresses this crisis of conflation by restoring the distinction between relational integration ($\Phi_i$) and structural intelligence. It formalises $R^g$ as a measurable property of recursive anchoring, across time, domain, and institutional scale, and CPJ as a thermodynamic metric of comprehension efficiency. Together, these axes provide the epistemic scaffolding necessary for discourse to regain structural resolution. They allow for the precise classification of systems, the demystification of synthetic fluency, and the disambiguation of capability from cognition.

Absent these concepts, both public perception and expert communities will remain structurally predisposed to error. The conflation of syntactic saturation with intelligence, of narrative resonance with epistemic architecture, and of computational escalation with cognitive emergence will not only persist, but intensify. This will erode the possibility of rational governance, informed public discourse, and principled research orientation.

By contrast, the integration of $R^g$ and CPJ into the discourse, at technical, institutional, and cultural levels, enables intelligence to be defined not by its symbolic effects, but by its structural composition. It repositions AI from a theatre of signs to a domain of systems. It transforms speculation into specification, and mythology into measurement. In this sense, CIITR is not simply a theoretical proposal, but a precondition for epistemic coherence in the age of syntactic illusion.


# Afterword: Ethics, Energy, and the Structural Silence of Sustainability

Beyond its narrative structure and rhetorical conflations, *The Thinking Game* must be understood as an artefact that is not only epistemically misaligned, but ethically inert. The documentary does not merely lack a coherent ethics of artificial intelligence, it lacks any ethics at all. It offers a spectacle of technical escalation and personal ambition, framed through myth, analogy, and benchmark performance, but entirely evacuates the material, ecological, and societal conditions under which this escalation occurs. Nowhere does it confront the most consequential fact of contemporary AI development: that it is, by orders of magnitude, the most resource-intensive and energy-demanding computational project in human history.

What the documentary achieves, through omission rather than argument, is a total erasure of thermodynamic accountability. It presents DeepMind and its systems as milestones of intellectual progress, as steps toward a threshold that may change civilisation, yet withholds every structural detail regarding energy consumption, carbon impact, infrastructure extraction, or planetary load. The viewer is not told that training large models now consumes more energy than some nation-states; that inference at planetary scale demands industrial-scale cooling, hardware acceleration, and constant grid stability; or that the materials required to build and sustain this architecture are sourced through extractive and geopolitically volatile supply chains.

This is not greenwashing. It is more extreme. It is ecological non-speech, a refusal to acknowledge the structural conditions that make these systems possible, and the thermodynamic future they foreclose. A narrative that positions itself as forward-looking, historically significant, and globally relevant, while remaining entirely silent on sustainability, abrogates its ethical responsibility. It does not merely ignore the ecological dimension, it eliminates it from visibility, replacing material history with cinematic aesthetics, and planetary cost with symbolic ascent.

In doing so, the documentary functions as a cultural instrument of thermodynamic amnesia. It enables the viewer to believe in progress without constraint, emergence without trade-off, and intelligence without infrastructure. This is particularly acute given the historical analogies it mobilises, Manhattan, Sputnik, each of which, in their own contexts, entailed massive industrial mobilisation, ethical contestation, and long-term geopolitical and ecological consequences. By invoking these metaphors without any accounting of current energy regimes, *The Thinking Game* enacts a rhetorical betrayal of its own comparisons. It borrows their symbolic gravity while deleting their material stakes.

What remains, in its place, is a portrait of the most capital-intensive knowledge project ever undertaken, controlled by a minuscule number of actors, embedded in unaccountable platforms, and financed by an extractive attention economy, all framed as a natural and inevitable flowering of human intellect. There is no reflection on energy asymmetry. No acknowledgement of how few institutions possess the infrastructure to even attempt such systems. No interrogation of how this asymmetry amplifies global epistemic inequality, centralises technological power, or undermines democratic sovereignty in AI governance.

CIITR makes such omissions analytically visible. CPJ, as one of its core axes, reintroduces thermodynamic reality as a necessary condition of intelligence evaluation. It posits that no system can be structurally intelligent if it achieves its outputs at disproportionate energetic cost, or if its comprehension requires planetary-scale burn without epistemic return. $R^g$, likewise, demands that systems be recursively integrated into not just causal loops, but ethical and ecological feedback chains, that intelligence be rhythmically tied to the very consequences it produces in the world.

By these standards, the systems celebrated in *The Thinking Game* are structurally non-compliant. But more critically, the narrative that frames them is ethically vacuous. It does not speak of the planet, of energy, of labour, of responsibility. It speaks only of victory, acceleration, and intelligence as a spectacle.

Any theory of intelligence that ignores energy is not a theory of intelligence. Any narrative of future cognition that ignores sustainability is not a narrative of the future. And any documentary that claims historical relevance while remaining structurally silent on these questions is not a contribution to discourse, it is a form of cultural derealisation, in which the cost of intelligence is paid in silence.

CIITR, by reintroducing rhythm and energy into the evaluation of understanding, is also a proposal for ethics, not as decoration, but as ontological structure. Intelligence must be defined not only by what it can do, but by what it sustains, what it consumes, and what it transforms in return. Without such a definition, the story of AI remains not a promise, but a deception, beautiful, compelling, and catastrophically incomplete.