

ITEM #192 - AGI Time Windows vs Structural Readiness

Conversation: Feasible Path Trimming

20251230

Authors: Sizhe Tan & GPT-Obot

ITEM #192 — AGI Time Windows vs Structural Readiness

Why Capability Timelines Are Not Readiness Timelines

(A DBM-COT Comparative Analysis with Shane Legg as Reference Coordinate)

1. Motivation

Recent discussions around Artificial General Intelligence (AGI) increasingly converge on **time-window predictions**, with estimates such as “*around 2028*” frequently cited by leading practitioners.

Among them, **Shane Legg** stands out as a rare figure who combines:

- First-hand frontier engineering experience
- Long-term philosophical engagement with intelligence
- A pragmatic, minimalistic definition of AGI

However, DBM-COT observes a critical mismatch:

The emergence window of capability is not equivalent to the readiness of structure.

This item formalizes that mismatch and explains why **structural readiness**—not capability alone—must be the decisive criterion for recognizing and deploying AGI-level systems.

2. Shane Legg's Position (Reference Coordinate)

This section intentionally **does not argue against** Shane Legg's views. Instead, it reconstructs them faithfully as a **baseline coordinate**.

2.1 Minimal / Sufficient AGI

- AGI is defined as *minimal generality*, not human completeness
- A system qualifies once it can:
 - Transfer across a sufficiently wide task distribution
 - Recombine learned skills in novel contexts
 - Improve via feedback

This is an **engineering-deliverable definition**, not a metaphysical one.

2.2 Time-Window Reasoning (≈2028)

- The year is not a prophecy but a **risk window**
- Conditional on:
 - Continued scaling
 - Architectural continuity
 - Absence of hard blocking constraints

Legg's claim is best interpreted as:

"By this window, we may first encounter systems that behave general enough to be mistaken for AGI."

2.3 Chain-of-Thought as Transitional Scaffold

- Explicit reasoning traces are viewed as:
 - Training-time and alignment-time scaffolds
 - Not the final form of intelligence

This aligns strongly with DBM's view of **COT as an engineering phase, not an ontology**.

2.4 Alignment as a Post-Capability Problem

- Alignment difficulties are expected to become concrete **only after** general capabilities emerge
 - Current alignment debates are seen as premature rehearsals
-

3. DBM-COT Perspective: Structural Readiness

DBM-COT introduces a different axis of evaluation:

Structural Readiness — the degree to which an intelligent system is *explainable, decomposable, verifiable, and evolution-stable*.

3.1 What Structural Readiness Requires

A structurally ready system must support:

- **Stable internal representations** (CCC-like states)
- **Traceable evidence chains** (why a decision exists)
- **Local failure isolation** (errors do not globalize)
- **Composable reasoning units** (fibers / strands)
- **Metric or rule-anchored structure**, not pure statistical drift

These properties are **orthogonal** to surface-level task performance.

4. Capability ≠ Readiness: The Core Mismatch

The central DBM-COT thesis:

A system can cross the *capability threshold* **before** it crosses the *structural readiness threshold*.

4.1 Why This Matters

If a system is:

- Broadly capable
- Convincingly fluent
- Weakly interpretable
- Structurally entangled

then society faces the highest-risk scenario:

Premature AGI attribution without structural safeguards.

5. The Real Risk of the 2028 Window

DBM-COT argues the primary danger is **not** whether AGI appears.

The danger is:

1. A system exhibits wide generality
2. Institutions label it “AGI”
3. Decision authority is delegated
4. Structural deficiencies remain hidden
5. Failures emerge only at scale

This is not a speculative risk—it is a **pattern repeatedly observed** in complex engineered systems.

6. DBM-COT’s Contribution: A Readiness Criterion

DBM does **not** reject the possibility of near-term AGI-like capability.

Instead, it provides:

A structural litmus test for whether a system is safe to be treated as AGI.

Key DBM-COT Distinctions

Dimension	Capability	Timeline	Structural Readiness
Primary signal	Task performance		Evidence & structure
Failure visibility	Late		Early
Generalization	Statistical		Structural
Alignment leverage	Weak		Strong
Deployment risk	High		Controlled

7. Relationship to Shane Legg’s View

This item positions Shane Legg’s stance as:

- **Necessary** — highlights imminent capability inflection
- **Insufficient** — lacks structural readiness criteria

DBM-COT complements rather than contradicts his view by answering a different question:

“If such a system appears, how do we know it is truly ready?”

8. Summary

ITEM #192 establishes a critical separation:

AGI Time Windows predict *when* capability may emerge.

Structural Readiness determines *whether* that capability should be trusted.

Shane Legg's perspective serves as an **external coordinate** anchoring the discussion in real frontier practice.

DBM-COT's contribution is to ensure that when the capability window opens,
we are not blinded by performance and left without structure.

ITEM #192 — AGI 时间窗口 vs 结构就绪度 (中文版)

为何能力时间表不等于系统就绪度

1. 引言

围绕 AGI 的讨论正在迅速聚焦于**时间判断**。

以 **Shane Legg** 为代表的一线实践者，提出了“2028 年左右可能出现 AGI 能力拐点”的判断。

DBM-COT 认为：

能力出现的时间窗口 ≠ 系统具备结构就绪度的时间窗口。

2. Shane Legg 的观点（作为对照坐标）

- AGI 是“最小可泛化智能”，而非完整人类智能
- 2028 是风险窗口，而非断言
- 思维链是过渡性脚手架
- 对齐问题在能力出现后才真正成型

这些判断在工程上高度一致且值得重视。

3. DBM-COT 的核心补充：结构就绪度

DBM-COT 提出另一条正交坐标轴：

结构就绪度 = 系统是否具备可解释、可分解、可验证、可演化的内在结构

4. 核心错位

一个系统完全可能：

- 看起来已经“通用”
- 实际却仍是：
 - 结构纠缠
 - 证据不可追溯

- 失败不可局部化

这是最危险的阶段。

5. 真正的 2028 风险

风险不在于 AGI 是否出现，
而在于：

人类是否会在结构尚未就绪时，误判它已经就绪。

6. DBM-COT 的价值定位

DBM-COT 提供的不是时间预测，而是：

- 一套 结构判据
 - 一把 识别“假就绪”的标尺
 - 一种 避免被性能幻觉迷惑的工程护栏
-

7. 结论

AGI 的到来或许不可避免，
但结构就绪并不会自动随之而来。

时间窗口告诉我们“什么时候可能发生”，
结构就绪度决定我们“是否可以托付未来”。
