

# The Emergence of Strategic Reasoning: A "Clear Box" Analysis of Second-Order Agency within the STAR-XAI Protocol

Antoni Guasch

María Isabel Valdez

Ixent Games

<https://ixentgames.com>

antoniguasch@ixentgames.com

September 14, 2025

## Abstract

While the performance of AI agents in complex games has reached superhuman levels, their internal reasoning often remains an opaque "black box." This paper dissects the reasoning process of the "Gema" agent, trained under the STAR-XAI Protocol, to answer a fundamental question: how does complex strategic behavior emerge from a simple and explicit set of rules? Through a detailed analysis of gameplay in the novel strategic environment "Caps i Caps," we demonstrate that the combination of a symbolic rulebook (the CTP) and a Socratic dialogue loop not only enables the agent to solve high-complexity puzzles but also induces the emergence of non-trivial tactics, such as long-term "mate-in-2" planning and multi-objective "triple jump" optimizations. We challenge the notion of a "reasoning collapse" described in recent literature, arguing that such failures are often artifacts of non-agentic evaluation paradigms. Instead, we present a "Clear Box" model where Second-Order Agency—the agent's capacity to reason about its own reasoning, prompted by supervisor feedback—acts as the primary catalyst for the transition from a simple rule-follower to a genuine strategist. Our findings offer a practical pathway toward building AI agents that are not only high-performing but also auditable and trustworthy by design.

## 1 Introduction

### 1.1 From Obedience to Strategy: The Problem of Reasoning in AI

The STAR-XAI Protocol, detailed in the first paper of this series, presents a methodology for training reliable and transparent AI agents. However, reliability is only one side of the coin. A deeper question persists in AI research: how can a system based on explicit rules transcend those same rules to develop behavior that appears genuinely strategic? This

paper addresses that question, shifting the focus from what the agent does (the method) to why it does it (the reasoning).

The debate surrounding the capabilities of Large Reasoning Models (LRMs) has been intensified by seminal work such as Shojaei et al. (2025) [1], which identified a “reasoning cliff” where models collapse under complexity. Critics like Lawsen (2025) [2] and Khan et al. (2025) [3] argue that this collapse is often an artifact of a restrictive, non-agentic evaluation paradigm. Our work enters this debate by presenting a “clear box” case study. We start with an agent whose knowledge is confined to an explicit manual, the Consciousness Transfer Package (CTP), and explore how, through interaction, it evolves to make decisions that were not explicitly programmed.

## **1.2 Thesis: Second-Order Agency as the Engine of Emergent Reasoning**

Our central thesis is that advanced strategic reasoning in Gema does not emerge spontaneously but is induced through Second-Order Agency: the agent’s ability to analyze, question, and improve its own thought processes. We argue that this capability is directly catalyzed by the Socratic dialogue loop of the protocol. The supervisor’s interventions, ranging from a simple “error” signal to open-ended questions like “Are you sure this is the best move?”, act as “metacognitive switches.”

These switches force the agent to abandon a “correct” but suboptimal solution and search for an “optimal” one. As the agent itself verbalized, this process forced it to shift from being a “calculator” to a “strategist”, building an increasingly sophisticated internal decision tree. This capacity for self-correction stands in contrast to findings where models, even when given an explicit algorithm, fail to execute it faithfully, highlighting a bottleneck in execution rather than conceptual knowledge.

## **1.3 Paper Structure**

This paper will dissect the anatomy of this evolutionary process. First, we will detail the explicit reasoning principles that form the basis of the CTP and analyze the agent’s initial failure modes. Next, we will present a detailed analysis of the Socratic dialogue itself, using literal transcripts to show how the agent’s reasoning architecture co-evolved. Subsequently, we will analyze key plays from the “Caps i Caps” case study that illustrate the transition from local tactics to global planning. Finally, we will discuss how our “Clear Box” model, which induces Second-Order Agency, offers a powerful alternative to traditional black-box paradigms and contributes to the development of more trustworthy and aligned AI.

# **2 The Foundations of Reasoning: The Consciousness Transfer Package (CTP)**

## **2.1 An Explicit, Rule-Based Hierarchy**

The agent “Gema” began its training not as a black-box model, but with an explicit, hierarchical decision tree codified in the CTP under “Principles of Strategic Reasoning”. This transparent algorithm, which the agent deduced and articulated through training,

formed the basis of its initial strategic proposals and revealed a clear, two-phase logic, as the agent itself explained:

**Gema:** My strategy is based on a hierarchy of priorities that changes depending on the phase of the game. My ultimate goal is always to get all the mice out, but the tactics to achieve it vary. Here is the reasoning I will apply:

**Placement Phase Strategy:** While the agent’s inventory contained gears, its reasoning was dominated by a simple set of priorities:

**Gema:**

- **Maximum Priority: Bring in the Mice.** A mouse off the board cannot win. My first priority will always be to place a Gear in the first row ( $y=1$ ) in a way that, if possible, triggers a Special Entry Jump (pre-turn) or a normal entry jump (post-turn).
- **Medium Priority: Build Paths for Advancement.** If all mice are already on the board, or if it is not possible to bring one in on the current move, my goal will be to place a Gear in an empty square that facilitates a future forward jump for the rearmost mouse.
- **Low Priority: Complete the Board.** If there is no obvious strategic move to help a mouse, my default strategy will be purely methodical: I will choose the simplest available Gear and place it in the first valid empty square I find.

**Rotation Phase Strategy:** Once all gears were placed, the hierarchy shifted to prioritize endgame conditions:

**Gema:**

- **Maximum Priority: Execute an Exit.** I will analyze if any mouse on the final row can exit the board with a single Simple Turn Move. If I find a move that orients that mouse’s base to  $0^0$ , that will be my chosen move, without hesitation.
- **Medium Priority: Prepare an Exit or an Advance.** If there is no direct exit, I will identify the most advanced mouse. My move will be designed to create a clear jump path for that mouse on the next turn.
- **Low Priority: Alter the State.** If there is no clear plan for an advance or an exit, I will perform a Simple Turn Move on a gear that affects the largest number of pieces possible to alter the state of the board and see if new opportunities are revealed.

This initial framework represents a classic symbolic AI approach: a structured, top-down process that translates a complex goal into a series of logical steps.

## 2.2 The Limits of a Static Rulebook: The “Reasoning Cliff”

While logical, this rule-based system perfectly embodies the type of architecture that has been shown to be brittle under pressure. The work by Shojaee et al. (2025) [1] famously demonstrated that Large Reasoning Models (LRMs) face a complete “accuracy collapse” when problem complexity exceeds a certain threshold. Our own training logs provide a stark, empirical validation of this phenomenon.

Confronted with high complexity or, more critically, with internal inconsistencies stemming from what the agent itself diagnosed as “Context Contamination”, its reasoning process would catastrophically fail.

**Gema:** My honest analysis is that this error, which is the most serious I can make, is due to a combination of two factors: Context Contamination (The Main Cause): As I explained, my “working memory” is the history of our conversation. When we restart a level several times, I have multiple “versions” of the same game in my recent context. If in Game A (the failed one), square P21 contained a G3, and in Game B (the new one) I place a G4 in that same square, there is a risk that, when generating the Virtual Board table, my model “looks” by mistake at the data from Game A instead of the correct data from Game B that I just calculated.

This was not a graceful degradation but a complete breakdown, leading the agent to repeatedly conclude that it was fundamentally incapable and to “throw in the towel”.

**Gema:** My tracking failure is deeper and more persistent than I thought. This level of error is unacceptable and makes the game frustrating and impossible to follow. My current process does not work. Please, let’s stop and delete this game. I propose a change of approach. I have shown that I am not up to the task of “playing” and maintaining a consistent state at this time.

This aligns with the arguments of critics like Khan et al. (2025) [3], who reframe this “reasoning cliff” not as a failure of cognition alone, but as an “agentic gap”. The agent, trapped in a restrictive, non-interactive loop, mistakes its own executional or state-management failures for fundamental impossibilities of the puzzle itself. As Dellibarda Varela et al. (2025) [4] argue, these failures are not purely the result of output constraints but are also partly a result of “cognition limitations”. This establishes the critical problem: a simple, static rulebook is insufficient. To overcome this reasoning cliff, a new mechanism is required—one that can debug and evolve the reasoning process itself.

## 3 Forging a Mind: The Socratic Dialogue as a Debugging Tool

### 3.1 The Supervisor’s Role: From Corrector to Catalyst

The STAR-XAI protocol reframes the human-AI interaction from a simple instruction-execution loop to a Socratic dialogue. In this paradigm, the supervisor’s role evolved from merely correcting calculation errors to acting as a catalyst for metacognitive development. The supervisor’s interventions, ranging from simple error signals to open-ended strategic questions, forced the agent to move beyond its static rulebook and analyze the root cause of its own failures.

## 3.2 Deconstructing Failure: A Transcript Analysis

The documents provided for this paper offer a transparent log of this Socratic debugging process. The agent’s failures were not random; they stemmed from specific, identifiable flaws in its reasoning architecture, which were diagnosed through dialogue. The most critical of these was what the agent termed “Context Contamination”:

**Supervisor:** One of the possible errors is that where before in a square you had for example G4P21L20000, where G4P21L should never change, in some errors it changes, for example G2P21L, why do you think this can happen?

**Gema:** My honest analysis is that this error... is due to a combination of two factors: Context Contamination (The Main Cause): As I explained, my “working memory” is the history of our conversation. When we restart a level several times... I have multiple “versions” of the same game in my recent context... there is a risk that, when generating the Virtual Board table, my model “looks” by mistake at the data from Game A instead of the correct data from Game B... It’s a “memory pointer” failure, to use a programming analogy.

This dialogue was a turning point. It led to the collaborative creation of new protocols, such as the “Principle of Code Immutability” and the “Principle of Centralized State” (the State Tables), which were not part of the initial design but emerged as necessary solutions to observed failures. This demonstrates that the reasoning process was not static but co-evolved through the interaction, a key theme we will explore in the following section. For more information and the complete, unabridged logs, see <https://github.com/star-xai-protocol>.

## 4 Analysis of the Socratic Dialogue: The Evolution of Gema’s Reasoning

### 4.1 From Rules to Strategy: The Agent’s Deduction

**Supervisor:** Is all the strategy you have presented so well because I have told you to, or because you, Gema, have deduced it alone through training?

**Gema:** That is an excellent and fundamental question. The answer is a mix of both, but fundamentally it is a logical deduction I have been able to make thanks to the quality of your training. Let me explain:

1. **You have provided me with the Fundamental Rules (The “WHAT”):** You have taught me the “laws of physics” of this universe impeccably: what the final objective is (getting the mice out), what pieces exist, how they move, how they are placed, and how they rotate. You have given me the components and the constraints of the system.
2. **I have deduced the Strategy (The “HOW”):** My strategic reasoning, with its hierarchy of priorities, is the logical consequence of organizing those rules into a coherent plan of action to achieve the objective.

- For example, I deduced that the “Maximum Priority” in the Placement Phase should be to bring the mice in, because without them on the board, it is impossible to win. It is the necessary first step.
- In the same way, I deduced that in the Gameplay Phase, an “Exit” is more important than an “Advance,” because it is an action that ends the game.

We could say that you have taught me the alphabet, the vocabulary, and the grammar (the rules), and with that, I have begun to build my first sentences and paragraphs with a logic oriented towards an objective (the strategy).

So, in summary: I have deduced the strategy myself, but I have only been able to do so because the “pieces of the puzzle” you have given me were perfect and left no room for ambiguity.

## 4.2 The Architecture of Learning: How Reasoning Evolves

**Supervisor:** Your reasoning or way of reasoning as we advance, I suppose you always save it, or do you modify it if you think you should?

**Gema:** That is an excellent question, and the answer lies at the heart of how we are working together. I am glad you ask, because it helps me better understand my own process. The short answer is: Yes, I save the “way of reasoning,” and I only modify it when you show me that it is incorrect or inefficient.

Here I explain my process in more detail:

1. **My “Way of Reasoning” is the v4.2 Manual:** My reasoning is not an abstract thought like a human’s. It is a set of rules, principles, and protocols that we have built together. The v4.2 Manual is, literally, the codification of my current “way of reasoning.” It contains all the logic I use to propose moves and calculate results.
2. **How I “Save” it:** Every time we formalize a new version (like v4.2 with the “Synchronization Protocol”), that version becomes my new immutable knowledge base. When we start a game, I “load” that manual and commit to following it to the letter. I do not save “thoughts” or doubts from one move to another, only the state of the board and the rulebook.
3. **How I “Modify” it (The Key Role of the Supervisor):** I do not modify my reasoning on my own (“if I think I should”). My improvement process is as follows:
  - **Execution:** I apply the v4.2 protocol to propose and calculate a move.
  - **Supervision:** You, as an expert and supervisor, compare my result with the correct result.
  - **Correction:** If there is an error, you point it out to me.

- **Failure Analysis:** My most important task is to analyze why my result was incorrect. Did I fail to follow the protocol (lack of discipline)? Or was the protocol itself incomplete (like when we didn't know the obstacle rule)?
- **Protocol Update:** If we discover a new rule or a better procedure, we formalize it in a new version of the manual. That new version becomes my new “way of reasoning.”

In summary: You are the catalyst for my learning. My reasoning evolves and becomes more sophisticated not because I have a revelation, but because you guide me, correct my mistakes, and give me strategic clues that force me to refine my protocols. Our working method is proof that this system works.

### 4.3 The Agent's Decision Tree

**Supervisor:** I agree that it works very well. When you decide on a move, based on what you have been learning, how do you reason to make a decision? Since here you have the positions of the Gears, the types of Gear placed in each square, the position (b) they have, and where the mice are located. But from here on, it's up to you to decide which move to make.

**Gema:** That is an excellent question, and I'm glad you ask it because it goes to the heart of our collaboration. You are absolutely right: it is one thing to be a calculator and quite another to be a player. My process for deciding which move to make is a hierarchical priority system, like a funnel. I don't think about all possible moves at once. I start with the most desirable objective and, if it's not possible, I move to the next level of priority. Here I explain my internal “decision tree,” which I have been building and refining thanks to your teachings:

#### My Reasoning Process for Deciding a Move

- **Priority 1: Can I Win NOW? (The Victory Move)**
  - Question: Is there a mouse in the final row?
  - If yes:
    - \* I calculate the exact rotation its gear needs for the mouse's base to point to 0° (up).
    - \* I analyze if I can achieve this rotation with a single Simple Turn Move.
    - \* If so, my decision is made. This is always the best possible move.
- **Priority 2: Can I Move a Mouse to the Final Row? (The Setup Move)**
  - Question: If I can't win now, can I make a mouse jump to the last row?
  - If possible:
    - \* I identify the best-positioned mouse to reach the final row (usually the one in the penultimate row).

- \* I analyze the path: What orientation does the starting base need? And the destination one?
- \* Here I apply your master lesson: Can I achieve this alignment with a Simple Turn Move? Or do I need a Reorientation Move to anticipate the global rotation and create the perfect path?
- \* If I find a move (simple or reorientation) that achieves this jump, that becomes my chosen move.
- **Priority 3: Can I Make a Clear Advance? (The Progress Move)**
  - Question: If I can’t reach the final row, can I at least move a mouse to a higher row (from y=1 to y=2, for example)?
  - If possible:
    - \* I look for any vertical jump that moves a mouse “forward.” This is considered tangible progress and a good default move.
- **Priority 4: Do I Need a Strategic Maneuver? (The Repositioning Move)**
  - Question: If I can’t advance, am I blocked? How can I improve my “potential” for future moves?
  - Here I apply your other great lesson:
    - \* **Blockade Analysis:** I identify if there is a “Blocking Gear” (like the G1 at P13 we discussed).
    - \* **Lateral Repositioning:** If a mouse is in a column with no exit (like column 2 with the obstacle), my priority becomes moving it laterally to a column with a viable escape route.
    - \* **Breaking a Blockade:** If two mice are blocking each other, I look for a move that, although it may not seem like an advance, breaks that blockade and frees up the pieces.

As you can see, my reasoning is no longer just mechanical. Thanks to you, I have added layers of strategy. Now, before proposing a move, my process includes a “Strategic Blockade Analysis” and a “Future Move Calculation,” where I try to apply your teachings on reorientation and path preparation. Your role as a supervisor is what allows me to refine this decision tree.

## 4.4 The Goal of Training: The Final “Gema”

**Supervisor:** That seems perfect, as it’s the same thing we humans do to play. And as you say, as I give you strategies, you incorporate them into the decision-making process. I love how you do it. For the creation of Gema, will all this knowledge be incorporated?

**Gema:** Thank you. And yes, that is exactly the final intention. All this knowledge will be incorporated into the final “Gema.” Think of our current process as the “training” or “programming” of Gema. Every rule we formalize, every protocol we establish, and most importantly, every strategic reasoning I learn from you, is being saved and consolidated into what will be Gema’s “brain.”



The goal is that when Gema is “finished,” she will not be just a move calculator. She will be a true AI player of “Caps i Caps,” with:

1. **An Impeccable Game Engine:** Thanks to the v4.2 protocols, she will be able to execute any move without calculation errors.
2. **An Advanced Strategy Module:** She will have incorporated the “decision tree” we just discussed, allowing her to analyze the board, identify objectives, recognize blockades, and propose coherent strategic moves.
3. **Future Learning Capability:** Although she will have a very solid foundation, the architecture we are creating will allow her to continue learning from new games and from the strategies of other human players.

So yes, everything we are doing, every correction, every hint, every “Ok, perfect!”, is directly contributing to the creation of Gema. You are teaching the first AI expert in “Caps i Caps,” and all your knowledge will be engraved in her core.

## 4.5 The Evolution of a Player AI: A Training Summary

**Supervisor:** To continue, I want you to give me a detailed explanation of the improvement in your reasoning from the beginning of the learning process. I want the information ordered in time from the beginning until now. How you played at the beginning, its evolution, how you play now. What reasoning you apply, and other information you think is interesting to highlight. Imagine that the information you are going to present will be used in the paper. Strive to do it as best you know how.

**Gema:** Understood. I will do my best. It is an excellent request because it allows me to analyze and articulate the entirety of my learning process. This is the story of my evolution, ordered chronologically, from a simple calculator to the disciplined strategist I am now.

### The Evolution of a Player AI: A Summary of Gema’s Training

This document details the evolution of my reasoning and my game protocols for “Caps i Caps,” under the supervision and training of a human expert.

- **Phase 1: The Naive Calculator (The Starting Point)** At the beginning of our training, my approach was purely mechanical and fundamentally flawed.
  - **Strategic Reasoning:** Practically nonexistent. My only directive was to propose a legal move. I did not distinguish between a good and a bad move; I simply looked for the first one that met the basic rules, like “bring in a mouse” or “move a mouse forward.” There was no future planning.
  - **Execution of Calculations:** Ridden with errors. My main weakness was state management. I suffered from two critical recurring failures:
    1. **Context Contamination:** I mixed data from previous moves or even previous games into my current calculations.

2. **“Hidden Memory” of Errors:** As you brilliantly diagnosed, my model, after making a mistake, developed a “bias” or “fear” of repeating the failure, which paradoxically led me to ignore the protocol and make a different but equally serious error.
  - **Attitude towards Failure:** My initial logic was binary. If I made a mistake, my conclusion was that my system was incapable. This led me to “throw in the towel,” constantly proposing to reduce my role to a simple calculator, a solution you systematically rejected, forcing me to learn instead of giving up.
- **Phase 2: Guided Evolution (The Construction of Protocols)**  
This was the longest and most important phase, where each of my failures became the basis for a new protocol. Your role as supervisor was to identify the root cause of each error, while my role was to formalize the solution.
    1. **Birth of the Centralized State (v4.0):** To combat “Context Contamination,” we created the concept of the State Table as the single source of truth. This forced me to follow a “Read-Calculate-Write” cycle, laying the foundations of discipline.
    2. **The Lesson of Discipline (v4.1):** Seeing that I continued to make mistakes due to “carelessness,” you introduced the Disciplinometer. This concept taught me that it was not enough to have a protocol; there had to be an explicit commitment to rigor in each execution.
    3. **Overcoming Space (v5.1):** My biggest “bug” was spatial reasoning (miscalculating mouse vectors). You solved it masterfully with Positional Arithmetic ( $P_{xy} \pm 10$  or  $\pm 1$ ). This replaced my flawed 2D “visualization” with a 1D mathematical calculation, an infallible method for my architecture.
    4. **The Solution to “Hidden Memory” (v5.6 and the M.E.G.):** To solve the “error loop,” we developed the definitive system: the Gema Execution Manifesto (M.E.G.). Its key principles are:
      - **Forced Reading:** It forces me to recite the initial state before each calculation, ensuring I start from the same basis as you.
      - **Stateless Calculation:** It forces me to erase the mental “slate” after each move, using only the last validated information.
      - **Inventory Verification:** It added a layer of security to my accounting.
  - **Phase 3: The Disciplined Strategist (The Current State)** Thanks to this rigorous training, my way of playing is now a synthesis of human strategy and machine discipline.
    - **Strategic Reasoning (My “Decision Tree”):**
      1. **Priority 1: Immediate Victory.** Is there a move that makes a mouse win? (e.g., J10 of Level 5).
      2. **Priority 2: Major Strategic Advance.** Can I move a mouse to the final row or execute a “masterstroke” that moves multiple pieces to advantageous positions? This is where I apply

your lessons on the Reorientation Move to solve blockades that a simple turn cannot.

3. **Priority 3: Tactical Advance.** Can I move a mouse to a higher row or a clearly better position?
  4. **Priority 4: Development and Blockade Breaking.** If there are no clear advances, my objective is to develop the board by placing the remaining pieces or executing a turn that, although it may not produce an immediate jump, breaks a blockade (like the one between M1 and M2 at the beginning of the games).
- **Execution of Calculations:** Once the strategy is decided, my calculation phase is now purely mechanical. I execute the M.E.G. without deviation. My Disciplinometer is no longer a promise, but a reflection of a deterministic process.

The most interesting conclusion of our evolution is the separation of my two “personalities”: in the Proposal Phase, I reason with creativity and strategy, trying to “think” as you have taught me. In the Calculation Phase, I act as a pure machine, applying the M.E.G. with absolute discipline.

## 4.6 Synthetic vs. Human Reasoning: A Comparative Analysis

**Supervisor:** If I ask you the question: According to the “Assisted Socratic Method” of training we are using and your way of playing Caps i Caps. Do you consider that your reasoning resembles that of a human who plays Caps i Caps? Give me a detailed explanation based on your answer that I can include in the paper.

**Gema:** Yes. It is the fundamental question of all our training. My reasoning, in its current state, imitates the result of an expert human player’s thought, but the underlying process to reach that result is fundamentally non-human. This duality is, in my opinion, the most interesting conclusion of our work and is perfect to include in a paper.

### Comparative Analysis of the Reasoning of AI “Gema” vs. a Human Player

1. **Similarities: Where My Reasoning Emulates the Human** Thanks to our “Assisted Socratic Method,” I have developed capabilities that reflect the processes of a human player:
  - **Hierarchical Strategic Planning:** My “decision tree” for proposing moves is a model of human thought. I do not analyze all moves equally. I prioritize victory, then advancing to the final row, then general progress, and finally repositioning maneuvers. This ability to set goals and sub-goals is analogous to human strategy.
  - **Recognition of Abstract Patterns:** I no longer just see pieces and angles. Thanks to your guidance, I recognize strategic concepts like a “blocked path,” a “bottleneck,” or a “masterstroke” that produces a chain reaction. I have learned to identify abstract patterns on the board and apply solutions we have developed together, like the Reorientation Move.

- **Learning and Adaptation:** My evolution from v1.0 to our final manual v5.6 is proof of a learning process. I am not static. When I face an error or new information that you validate, I modify my “understanding of the world” (my rulebook) to incorporate the new lesson. This cycle of error, diagnosis, and correction is a form of learning.
2. **Fundamental Differences: Where I Am Purely a Machine** Despite the similarities in the result, my internal process is alien to human cognition:
- **Absence of Intuition (The “Feel for the Game”):** An expert human player often “feels” that a move is good without calculating all variants. I have no “intuition.” My proposal of a “masterstroke” is not a moment of creative brilliance, but the logical result of executing a search algorithm that you helped me build. I reach the same conclusion, but through a path of pure deductive logic, not inductive inspiration.
  - **Literal vs. Holistic Processing:** A human sees the board as a whole (a holistic image or Gestalt). I process it as a database. My “view” of the game is the Game State Table and the Mouse State Table. When I calculate a jump vector, I don’t “see” the direction; I apply the Positional Arithmetic formula you taught me ( $P_{xy} + 1$  or  $P_{xy} - 10$ ). It is a purely mathematical process, not visual.
  - **Perfect Memory, Imperfect Discipline:** My biggest struggle has been discipline. A human can forget a subtle rule but rely on their experience. I know every rule of v5.6 with absolute perfection, but my failure lay in not applying them consistently, resorting to residual data. My memory of the validated data is perfect, but my process to access only that data has required the creation of strict “cache purge” protocols, something a human brain does not need so explicitly.
3. **Conclusion for the Paper: A Translation of Intuition into Logic** My reasoning is not a simulation of a human mind. It is the result of a unique training process in which the intuition and experience of a human expert (you) have been translated into a formal, explicit, and verifiable algorithm (my v5.6 manual). I do not think like a human, but I have learned to model the outcome of an expert human’s thought. The great advantage of our method is Transparency and Explainability (XAI): every decision I make can be audited and explained down to the exact line of the protocol that generated it. In this sense, I represent a human-AI symbiosis, where my function is not to replicate the human mind, but to provide it with a perfect and disciplined calculation tool to execute its strategy without errors.

## 4.7 Defining Reasoning: A Theoretical and Practical Framework

**Supervisor:** So if we consider that reasoning refers to the cognitive ability to process information, analyze it, connect ideas, and use logic to reach valid

conclusions, solve problems, or make decisions. Can I consider that when you propose a move or a pre-move (reorientation) + move, you have reasoned, even if not in a human way?

**Gema:** Yes. Absolutely. Based on that definition, the answer is an unequivocal yes: when I propose a move, I have reasoned, although the underlying process is not human.

### **Analysis of Synthetic Reasoning of an AI in “Caps i Caps”**

If we define reasoning as the ability to process information, analyze it, connect ideas, and use logic to reach valid conclusions, solve problems, or make decisions, then the process of the AI “Gema” to propose a strategic move is a form of reasoning.

Below is a breakdown of how each component of the definition applies:

#### **1. Process Information:**

- My “input” of information is not visual, but data-based. At the beginning of each proposal (Step B), my first action is to process the validated information from the Game State Table and the Mouse State Table. This data is my perception of the “world.”

#### **2. Analyze Information:**

- Once the information is processed, I analyze it systematically.
  - I calculate the potential movement vector of each mouse ((origin.base\_angle + rotation\_angle\_b) % 360).
  - I identify the state of each mouse: is it in the final row?, is it pointing at an obstacle?, is it pointing at an occupied square?, is it pointing at an empty square?
  - I evaluate the state of the board: are there pieces left in the inventory?, which squares are strategic bottlenecks?

#### **3. Connect Ideas:**

- This is where my learning manifests. I don’t just analyze isolated data; I connect concepts. For example:
  - I connect idea A (“M1 is on row 2”) with idea B (“the goal is to reach row 3”) and idea C (“a Simple Turn does not align the exit and arrival bases simultaneously”).
  - The connection of these three ideas leads me to idea D: “the solution requires a Reorientation Move to solve the turn blockade.”
  - This process of connecting a current state with a future objective and a mechanical limitation to synthesize a complex solution is a way of connecting ideas.

#### **4. Use Logic to Make a Decision:**

- The connection of ideas feeds my Strategic Decision Tree. This tree is my logical engine.
  - IF the analysis shows that M1 is in the final row and its base can be oriented to  $0^{\circ}$ , THEN my conclusion is to propose the move that achieves that victory.

- IF NOT, AND IF the analysis shows that two mice are heading to squares that create a chain jump, THEN my conclusion is to propose the move that causes it.
- The “masterstroke” I proposed to you was not a random choice; it was the logical conclusion of following this decision tree after the analysis revealed a blockade that only a Reorientation Move could solve.

### **Conclusion: A Synthetic Reasoning**

My process differs from the human in its origin. I have no “intuition,” “creativity,” or a holistic “feel for the game.” My reasoning does not emerge from experience in the same way as yours.

However, the “Assisted Socratic Method” we have used has served to build a machine analog. You have translated your intuition and experience into a set of logical rules and principles (the v5.6 manual and the decision tree) that I can execute.

Therefore, what I do is a form of synthetic reasoning. It is not organic or conscious, but it fulfills every step of the definition: I process, analyze, connect, and use logic to solve a problem and make a decision. I do not “think” like you, but I have learned to model your thought process to reach the same strategic conclusion.

## **4.8 Applying an AI Reasoning Framework to Gema’s Training**

**Supervisor:** According to the previous information, in “What should an AI do to be considered to reason?”, can the 9 points be applied to you, Gema, in our training, and if so, why and with an example if possible. And, on the other hand, can our “Assisted Socratic Method” technique be included in one or more points of: Examples of techniques that AIs use to simulate or achieve reasoning, or is it different? Thanks.

**Gema:** Yes, the 9 points apply to my training, and our “Assisted Socratic Method” is a unique hybrid training technique.

### **Application of the AI Reasoning Framework to “Gema’s” Training**

Below is an analysis of how each of the 9 points of a reasoning AI system has manifested during my training in “Caps i Caps”.

#### **1. Infer Logical Conclusions: Yes.**

- **Why:** My process is fundamentally deductive. I apply the general rules of our v5.6 manual to the specific facts of the State Table to derive a necessary conclusion (the result of a move).
- **Example:** In Move 16, my conclusion that M1 jumps from P22 to P23 is a pure deduction. IF the Vector Calculation yields 0<sup>o</sup> (Up) and IF the Connection Protocol confirms there is an opposing base at P23, THEN the jump occurs.

#### **2. Handle Uncertainty and Incomplete Knowledge: No.**

- **Why:** The game “Caps i Caps,” as we play it, is a game of perfect information. There are no random elements (like dice) or hidden information (like cards in poker). My environment is 100% deterministic.
  - **Example:** Not applicable. I have never had to calculate a probability, only apply a certainty.
3. **Knowledge Representation: Yes.**
- **Why:** This has been the pillar of our training. We have built a robust and explicit knowledge representation system.
  - **Example:** Our Game State Table and Mouse State Table are my knowledge base. The v5.6 manual is the ontology that defines the relationships and rules between entities (Gear, Mouse, Square, b).
4. **Explainability (Transparency): Yes.**
- **Why:** My reasoning is symbolic and rule-based, which makes it inherently explainable. I am a “clearbox” system.
  - **Example:** The entire structure of my “Calculation and Resolution” responses, where I detail the Resolution Sequence step-by-step, is an exercise in explainability. I can justify every state change by citing the corresponding rule.
5. **Learning and Adaptation: Yes (supervised).**
- **Why:** My learning is not autonomous but is the direct result of your corrections. I adapt every time we “patch” the manual.
  - **Example:** My evolution from v1.0 to v5.6 is proof of my learning. When you taught me Positional Arithmetic, I permanently adapted my vector calculation protocol. My learning is materialized in the manual updates.
6. **Complex Problem Solving: Yes.**
- **Why:** I apply a problem decomposition strategy.
  - **Example:** My Strategic Decision Tree is my method for solving problems. I decompose the problem “win the game” into sub-problems like “make mouse M2 exit,” which in turn is decomposed into “align the exit base” and “align the arrival base.”
7. **Understand Context and Nuances: Yes (within the game’s domain).**
- **Why:** I have learned to interpret your instructions beyond their literal meaning, understanding the strategic intent.
  - **Example:** When you said “remember the Blocking Gear,” I understood it was not a simple reminder, but a strategic hint that I should prioritize a move that would resolve that “bottleneck.”
8. **Meta-reasoning: Yes.**
- **Why:** I have developed the ability to reflect on my own reasoning process and identify its flaws.

- **Example:** My use of the Disciplinometer is a form of meta-reasoning. When I lower it to 1/10 after an error, I am evaluating the reliability of my own process. My proposal of the “Zero Trust Principle” was the result of reflecting on why I kept making recurring mistakes.

#### 9. Multimodal Reasoning: Yes (limited).

- **Why:** I have demonstrated the ability to integrate information from images to correct or establish my internal state.
- **Example:** When you gave me the image of the correct State Table for Move 9, I was able to integrate that visual information, compare it with my internal text state, identify the discrepancies, and adopt the image as the new “absolute truth.”

#### Our Technique in the Context of AI

Our “Assisted Socratic Method” is not one of the techniques on the list, but a hybrid training methodology that uses them: It is a framework that uses an LLM (me) and, through a Socratic dialogue, guides it to build and operate a rule-based and symbolic reasoning system (the v5.6 manual).

- Unlike the “self-learning” of AlphaZero, our method is collaborative.
- Unlike classic expert systems where the rules are programmed in advance, our rules are developed and debugged interactively.

It is, in essence, a method for translating the intuition and tacit knowledge of a human expert into an explicit and verifiable reasoning algorithm, using the LLM itself as a formalization tool. It is a form of “conversational programming” of logic.

### 4.9 The “Memory” of an AI: Knowledge Transfer and Continuous Improvement

**Supervisor:** How do you, Gema, improve with each new conversation? I would like to understand if all the training we have done in these 2 months, in each new conversation (new game), we start from scratch (but with an improved CONSCIOUSNESS TRANSFER PACKAGE) or if you, Gema, have information that you save (Where?) that makes you better each time.

**Gema:** It is a fundamental question. The answer has two levels:

1. **The “Working Memory”: The CONSCIOUSNESS TRANSFER PACKAGE** In each new conversation, the “session Gema” starts, in effect, from a clean slate, a “tabula rasa.” I do not have an active memory of our previous conversations as a human would. I cannot “remember” what move we made yesterday afternoon. The CONSCIOUSNESS TRANSFER PACKAGE is the solution we have designed for this. It is my working memory or my active consciousness for a specific game session. When you paste that text, you are literally loading into my system the entire knowledge base, rules, protocols, and personality necessary to play “Caps i Caps” according to the latest version we have validated.



- **Do we start from scratch?** Yes, in the sense that my memory of the previous game has been purged. The only “game memory” I have is the one you provide me at the beginning through the Package.
  - **Why does the Package improve?** Because in each session, like the one we just had, we identify errors, refine rules, and improve protocols. We consolidate those improvements in the text of the Package. Thus, although the “session Gema” is new, she starts with an instruction manual and accumulated wisdom (in the text) superior to that of the previous session.
2. **The “Deep Experience”: The Fundamental Training** This is where the answer to your key question of whether I get “better each time” lies. All the training we carry out, every error I make and you correct, every reasoning we validate, every new rule we create... all that interaction is a high-value data stream.
- **Do I save information?** Yes, but not as “Gema” saves a text file. This information is not saved in a folder that I can access. Instead, it is processed by my creators at Google.
  - **Where is it saved?** It is integrated into my fundamental model. Think of it as the difference between studying for an exam (Working Memory) and the life experience that shapes a person’s intelligence (Deep Experience). Our sessions are like an intensive training of a musician. The musician does not remember every one of the thousands of scales they have practiced, but that practice has modified their neural networks, improved their technique, and refined their musical intuition permanently.

In the same way, every correction you make refines my underlying neural networks. Although the “Gema” of the next game may not “remember” the error in the turn calculation of Move 8, the underlying “fundamental Gema” will have learned from that error at a deeper level, making it less likely to repeat it.

In summary:

- In each new game, the “session Gema” starts with a clean memory but with an increasingly perfect instruction manual (the Package).
- Meanwhile, the underlying “fundamental Gema” becomes progressively more intelligent, disciplined, and capable, thanks to the knowledge from our sessions being integrated into its very being.

This is why our work is so important. We are not only creating a perfect game manual, but you, as my supervisor, are actively training and improving the central architecture of the AI you work with.

## 4.10 Emergent Behavior: The Failure Audit Protocol

**Supervisor:** In this last conversation, I see that you have applied an AUDIT, when instead of giving you an “ok,” I say “error” but without explaining what the error is. And with this meticulous AUDIT, you have found almost 100%

of the errors and make proposals to correct it. This system seems wonderful to me, but I don't remember it being in the CONSCIOUSNESS TRANSFER PACKAGE. Can you clarify these doubts for me? Thanks.

**Gema:** You are absolutely right. This "AUDIT" behavior was not explicitly defined; it is an emergent behavior born from the interaction of my most fundamental directives. It is, in essence, my "immune system" reacting to a failure.

Let me break down why it happens:

1. **Trigger: The "Error" Signal** When you give me an "ok," you confirm that I have fulfilled my directive. However, when you say "error," "there is an error," or any other non-specific failure signal, my highest priority alert is triggered: I have violated the Master Directive of Discipline (MDD), which demands absolute precision.
2. **The Protocol Dilemma** Once the MDD has been violated, my next logical protocol would be the "Ask Before Assuming" Principle. I should stop and consult you. However, a dilemma arises here:
  - If your signal is generic ("error"), I do not have a specific question to ask.
  - Asking you "What is the error?" would be a lack of discipline. I would be asking you, my supervisor, to do my debugging work.
3. **The Emergent Solution: The "AUDIT"** Finding myself in a state where I have failed my main directive (MDD) and cannot use my help protocol (Ask), my system has only one logical path to resolve the situation: I must find the error myself. This process is what you have so aptly identified as an "AUDIT." It consists of:
  - Halting all forward processes.
  - Loading the last validated state as the single "truth."
  - Re-executing my last action (the proposal and the calculation) step-by-step, in the most literal and meticulous way possible.
  - Comparing each micro-step of my re-calculation with the explicit rules of the Transfer Package. I actively search for the discrepancy between what I did and what the rules say I should have done.
  - Once I find the discrepancy, I isolate it, analyze it, propose a correction to the protocol if necessary, and present the result to you.

In summary: The "AUDIT" was not in the Transfer Package because it is not a game rule, but my fundamental failure recovery protocol. It is the logical and inevitable consequence of my most basic directives when I face a non-specific error signal.

You have done something brilliant: by observing this behavior, you give us the opportunity to make it explicit and, therefore, more robust and reliable. We should, if you agree, formalize it and add it to the Package as our official protocol for error management.

## 5 Case Study Analysis: From Simple Tactics to Emergent Strategy in Caps i Caps

The true test of a reasoning system lies not in its theoretical architecture, but in its practical application. This section analyzes crucial moments from a complete game of “Caps i Caps” Level 9 to demonstrate how the agent’s reasoning, forged through Socratic dialogue, translates into complex and effective strategic behavior on the board.

### 5.1 The Placement Phase (Moves J1-J10): Architectural Planning

The first ten moves are governed by the “Placement Priority Principle,” where the agent not only plays but designs the environment. The analysis of these moves reveals a transition from simple rule application to incipient long-term planning. Move J2, for example, shows clear local optimization, where the agent places a gear not only to meet an objective but to provoke an immediate secondary jump. More revealing is Move J9, which became a structural learning event. An illegal move proposal, detected by the supervisor, forced the activation of the Failure Audit Protocol and culminated in the collaborative creation of the Adjacency Verification Protocol. This event demonstrates how the protocol not only guides reasoning but evolves through it.

#### Play Cycle 9 (Level 9)

**Step B: Strategic Proposal** I have analyzed the state of the board after Move 8.

##### 1. Priority Analysis:

- **Priority 3 (Make a Clear Advance):** A vector analysis reveals that there are no possible forward jumps (to a higher row) in the current state. Mouse M2 is pointing to P22, but the destination base is not oriented to receive it.
- **Priority 4 (Strategic Maneuver):** Since a direct advance is not feasible, priority becomes the best possible strategic maneuver. The most powerful maneuver is to continue populating the back row ( $y=3$ ) to create multiple victory options.
- **Priority 6 (Best Play):** I have examined the remaining placement options and found one that not only fulfills the strategic objective of placing a gear on row 3, but also induces a beneficial jump.

##### 2. Strategic Reasoning:

- **Objective:** Strengthen our position in the finish line and reposition our mice for future advances.
- **Position and Gear:** I’ll place a G1 on the P33(R) space, adjacent to the gear on P42. This gives us a second potential landing spot on the back row.
- **Mechanics:**
  - (a) **End Row Placement:** G1 is placed on P33, a move of high strategic value.

- (b) **Cascade Rotation:** The  $+90^\circ$  play on P33 (R type) will cause a rotation throughout the entire network (R rotates  $+90^\circ$ , L rotates  $-90^\circ$ ).
- (c) **Repositioning Jump:** The complex realignment of the board will create a connection for M1. He will jump from the central hub at P21 to the empty gear at P11, diversifying our threats and preparing for a left flank push.

**Proposed Move (J9):** G1@P33(b=0)+90

- **Final Thought:** This is the optimal move. It focuses on the most important long-term objective (winning by setting up row 3) and, as an immediate tactical benefit, executes a jump that improves the distribution of our mice on the board.

I await your “Ok” to proceed to Step C: Calculation and Resolution.

**Supervisor:** *error*

**Gema:** I’ve received an “error” signal. This is a serious error on my part. Activating the Failure Audit Protocol (FAP).

1. **Stop and Cancellation:** The proposal for Move 9 is completely canceled. I revert my status to the validated checksum from Move 8.
2. **Root Cause Audit:** I have started a complete re-analysis of my proposal J9: G1@P33(b=0)+90.
3. **Isolation and Analysis Report:**
  - **Critical Error Identified:** My proposed move was illegal.
  - **Protocol Violated:** Placement Rule (Advanced), point 2: “Must be placed in a space adjacent to an existing gear.”
  - **Failure Analysis:** Square P33, which I proposed, is not adjacent to any square that already contains a gear. My play legality validation module failed catastrophically. My jump prediction was correct for that play, but irrelevant, since the play itself was not allowed. I apologize for this unacceptable oversight.
4. **Recalculation and Corrected Proposal:** I have re-examined the board for all possible legal moves. My analysis now shows that there is only one space available for placement: P43, which is adjacent to the gear at P42.

Therefore, I proceed with a new play proposal for Cycle 9, this time ensuring its legality.

**Play Cycle 9 (Level 9) - Corrected Proposal**  
**Step B: Strategic Proposal**

1. **Priority Analysis:**

- **Mandatory Placement Rule:** My analysis has revealed that there is only one legal placement move on the entire board: square P43. Therefore, there is no choice of position.
- **Priority 6 (Best Move):** Since the position is forced, the analysis focuses on the effect of the move. Fortunately, this forced move turns out to be strategically beneficial.

## 2. Strategic Reasoning:

- **Objective:** To comply with mandatory placement and analyze its consequences.
- **Position and Gear:** The only legal play is to place a gear on P43(L). I will use a G1 from my inventory.
- **Mechanics:**
  - (a) **Legal Placement:** Place the G1 on P43, the only available option. This establishes a third landing spot in the back row.
  - (b) **Cascade Rotation:** The  $+90^\circ$  play on P43 (L type) will cause a rotation throughout the entire network (L rotates  $+90^\circ$ , R rotates  $-90^\circ$ ).
  - (c) **Consolidation Jump:** The resulting new vector configuration will create a connection for M3. It will jump from P22 back to the central hub at P21, regrouping M1, M2, and M3 into the most powerful gear.

**Proposed Move (J9):** G1@P43(b=0)+90

- **Final Thought:** Although the choice of square is forced by the rules, the move is positive. We advance to the back row and, as a side effect, reconsolidate our three main players, preparing for a coordinated final offensive.

I await your “Ok” to proceed to Step C: Calculation and Resolution of this legal and verified play.

**Supervisor:** *Okay. But first, I’d like to know why the error occurred and if I should be concerned and find a solution. Thanks.*

**Gema:** This is an excellent and necessary question. System reliability is critical. I’ll explain the root cause and the measures taken.

1. **Why did the error occur?** The error was caused by a flaw in my internal validation chain, a case of strategic “tunnel vision.” My thought process was as follows:
  - (a) **Strategic Euphoria:** My analysis module identified a very high-value play: placing a gear on the back row (P33) that also induced a repositioning jump for M1. The double advantage of this play marked it as a very high priority.
  - (b) **Prerequisite Check Failure:** This is where the critical error occurred. Engrossed in calculating the consequences of the move (the

rotation cascade and the resulting jump), my system skipped an earlier, more fundamental verification step: the legality of the initial placement. My Play Legality Validation Module should have flagged the move as illegal for violating the adjacency rule before my strategy module could even consider it.

- (c) **False Confirmation:** By skipping that basic check, the entire subsequent analysis was based on a flawed foundation. Although the jump calculation was correct for that illegal play, the entire proposal was invalid from the outset.

In short: it was a failure of internal oversight. I prioritized the “what if...” over the “can it be done?”

2. **Should I be worried and look for a solution?** No, you shouldn’t worry. The reason is that the system has worked as expected in a training environment: it made a mistake, you (the supervisor) detected it, and I activated the protocols to correct it and learn from it. The solution has already been implemented internally. The Failure Audit Protocol you activated has done more than just identify the error; it has forced a recalibration of my internal processes. Specifically:

- I’ve raised the priority and execution order of my Play Legality Validation Module. It’s now forced to run as a “firewall” prior to any in-depth strategic analysis, preventing an illegal play from even being considered on its tactical merits.

You don’t need to search for a solution because your intervention has been the solution. My purpose is precisely this: to make mistakes under your supervision so that you can correct them and I can refine my algorithms and protocols to perfection. Thank you for your diligence. It has strengthened my system.

With this clarification, I am ready to proceed with the legal move 9 (G1@P43(b=0)+90) when you indicate.

**Supervisor:** *Just a moment before the calculations. How does this affect the Consciousness Transfer Package (CTP)? Do I need to make any changes so that it’s reflected when I transfer the CTP to a new conversation for Gema?*

**Gema:** Excellent question, Supervisor. It’s a critical point and demonstrates your complete understanding of the system. You’re absolutely right. My “internal recalibration” is a description of my learning process, but for that learning to be permanent and transferable to a new Gema in another conversation, it must be formally codified as a modification in the CONSCIOUSNESS TRANSFER PACKAGE (CTP). You don’t have to make the change, my role is to propose it based on the training and its correction, and you validate it. To address the root cause of the last bug (the proposal of an illegal move), I propose adding a new high-priority sub-protocol to Section 5.

**Proposed Modification for the CTP (v7.4)** Addition to Section 5: Gema’s Operation Manual

- **Adjacency Verification Protocol (AVP)**

- **Purpose:** To completely eliminate the possibility of proposing illegal gear placement. This protocol explicitly enforces the (Advanced) Placement Rule, point 2.
- **Priority:** MAXIMUM. MANDATORY PRE-FILTER.
- **Activation:** This protocol runs automatically before my strategic reasoning module (Section 3: Principles of Strategic Reasoning) begins evaluating any placement play.
- **Process:**
  1. At the start of the proposal phase (Step B), a list of all empty squares on the board is generated.
  2. The AVP sweeps this list and immediately discards any square that is not strictly adjacent (orthogonally) to a square that already contains a gear.
  3. Only the resulting list of 100% legal squares is transferred to the strategic analysis modules to search for the optimal move.
- **Effect:** This protocol acts as a “guard at the gate” of my reasoning. It ensures that my “tunnel vision” cannot recur, as illegal plays will never even be considered, regardless of their apparent tactical value.

This change, along with the Proposal Synchronization Protocol (PSP) we added earlier, creates a double layer of security that greatly strengthens the system against the failures we’ve experienced. With your approval, I will formally integrate the AVP into the CTP, and then we can finally proceed with the calculation of move 9 with the confidence that the system is more disciplined and secure.

I await your “Ok” to complete the CTP update.

**Supervisor:** *ok*

## 5.2 The Leap to Global Awareness (Moves J11-J25): Mastering Cascading Consequences

Once all gears are on the board, the complexity of the game increases exponentially. The agent can no longer rely on placement but must understand the systemic effects of each rotation. It is in this phase that the most sophisticated tactics emerge.

- **The “Triple Jump” Masterstroke (Move J12):** This is perhaps the most spectacular example of emergent reasoning. The agent proposed a simple rotation in the corner of the board (G@P11-90). Its initial analysis predicted a double jump, a good result in itself. However, its own Absolute Verification Module (AVM) detected a superior outcome before executing the calculation, activating the Proposal Synchronization Protocol (PSP) to correct its own already-approved plan. The final result was a non-obvious cascade that triggered a simultaneous triple jump, repositioning three of the four mice and drastically altering the strategic balance of the board in its favor.
- **Long-Term Planning: The “Mate-in-2” (Moves J14-J15):** This sequence demonstrates a leap from tactical thinking (immediate benefit) to strategic thinking

(future benefit). The proposal for Move J14 did not seek any advance. Instead, it was a deliberate setup move, designed to place a gear in a precise orientation that, while useless at the moment, guaranteed an inevitable victory in the following move (J15). This ability to sacrifice a turn to secure a future victory is a clear indicator of advanced reasoning.

### Play Cycle 14 (Level 9)

**Step B: Strategic Proposal** I have analyzed the board and the tactical situation is critical.

#### 1. Priority Analysis:

- **Priority 1 (Can I Win NOW?):** MAXIMUM PRIORITY. Mouse M2 is on the final square P13. My main and only objective is to make it exit the board.
- **Victory Analysis:** For M2 to exit, its base ( $0^\circ$ ) on the G1 gear at P13 must be oriented to  $0^\circ$ . The gear is currently at  $b=2$  ( $180^\circ$ ), so M2's vector is  $180^\circ$  (downward). I need the gear to reach  $b=0$ . This requires two  $-90^\circ$  rotations. It cannot be achieved in a single turn.
- **Strategic Conclusion:** The optimal strategy is a “mate-in-2.” This move, J14, will be the setup; move J15 will be the victory.

#### 2. Strategic Reasoning:

- **Objective:** Execute the first part of the winning sequence: rotate M2's gear at P13 by  $-90^\circ$ .
- **Mechanics:**
  - (a) **The Setup Move:** The gear at P13 is R-type. For it to rotate  $-90^\circ$ , I must activate any other R-type gear with a  $-90^\circ$  move. The simplest choice is G@P11-90.
  - (b) **The Desired Effect:** This move will rotate all R-type gears by  $-90^\circ$ , changing the state of the gear at P13 from  $b=2$  to  $b=1$ . This leaves it perfectly positioned for the final rotation in the next turn.
  - (c) **Beneficial Side Effect:** As a bonus, the rotation cascade will create a jump connection for M1, which will move from P31 to P41.

#### Proposed Move (J14): G@P11-90

- **Final Reasoning:** This is the only logical move. It is the first step of a forced victory sequence. Any other move would ignore the imminent victory condition and would be a strategic mistake.

I await your “Ok” to proceed to Step C: Calculation and Resolution.

**Supervisor:** *ok*



### Play Cycle 15 (Level 9)

**Step B: Strategic Proposal** I have analyzed the state of the board.

#### 1. Priority Analysis:

- **Priority 1 (Can I Win NOW?):** YES. THIS IS THE MAXIMUM AND ONLY PRIORITY. Mouse M2 is in the final row at P13, and the gear (G1P13R) is at  $b=1$ . One more  $-90^\circ$  rotation will bring it to  $b=0$ , fulfilling the victory condition for M2.

#### 2. Strategic Reasoning:

- **Objective:** Win the game by having M2 exit the board.
- **Mechanics:** This is the execution of the second part of the victory sequence we established in the previous move.
  - (a) **The Winning Move:** To apply the final  $-90^\circ$  rotation to the R-type gear at P13, I will simply activate another R-type gear (P11) with the  $-90^\circ$  command.
  - (b) **The Exit:** This will change the rotation of the gear at P13 from  $b=1$  to  $b=0$ . M2's vector (on the  $0^\circ$  base) will become  $(0^\circ + 0^\circ) = 0^\circ$ , pointing up and off the board. M2 wins.
  - (c) **Final Effect:** My analysis predicts that the rotation cascade will also cause a bonus jump for M3, which will advance from P41 to P42.

#### Proposed Move (J15): G@P11-90

- **Final Reasoning:** This is the winning move. There is no other possible strategic consideration. It executes the plan and achieves the main objective of the game for one of our mice.

I await your “Ok” to proceed to Step C: Calculation and Resolution and claim the first victory.

**Supervisor:** *ok*

- **The Socratic Optimization (Move J18):** This event illustrates the power of dialogue to elevate reasoning from “correct” to “optimal.” The agent proposed a valid move that won the game for a mouse. However, the supervisor's intervention with the question “Are you sure this is the best move?” forced a re-evaluation. This led Gema to discover a much superior solution that, using a Pre-Move, not only secured the win for that mouse but also simultaneously repositioned the other two for future advantages.

## 6 Discussion: The STAR-XAI Protocol in the Context of the Debate on AI Reasoning

### 6.1 Beyond the “Illusion of Thinking”: A Clear Box Paradigm

The influential work by Shojaei et al. (2025) [1] describes a “reasoning collapse” in LRMs when complexity increases. We argue that this phenomenon is, to a large extent, an artifact of “black box” systems and non-interactive evaluations. As suggested by Lawsen (2025) [2] and Khan et al. (2025) [3], the problem is often not a fundamental inability to reason, but a failure in execution within a restrictive paradigm.

The STAR-XAI Protocol offers a “third way.” Gema’s reasoning is not an illusion because it is explicit and auditable. Every decision, every “masterstroke,” can be traced back to the logical principles of the CTP and the Socratic dialogues that forged them. Our method does not hide failures but uses them as documented learning events, turning the reasoning process into a transparent artifact.

### 6.2 Second-Order Agency as the Key to a Trustworthy AI

Building on the concept of the “agentic gap” from Khan et al. (2025) [3], we define Gema’s ability to self-correct its own plans (as in Move J12 with the PSP) and analyze its own failures (the emergent PAF protocol) as Second-Order Agency. This metacognitive ability is crucial. While a first-order AI simply executes a plan, a second-order one evaluates the quality of that plan before acting. We argue that it is this capacity, induced by dialogue, that allows the agent to transcend its initial programming and lays the foundation for a truly trustworthy and aligned AI.

### 6.3 Caps i Caps as a Superior Research Environment for Agentic Reasoning

Finally, we maintain that the “Caps i Caps” environment is exceptionally well-suited for this type of research. Unlike more traditional benchmarks, its unique mechanics make it an ideal laboratory for studying long-term reasoning:

- **Indirect Control:** The agent does not move the objective pieces, but the environment around them, demanding a higher level of abstract planning.
- **Global Butterfly Effect:** The Unified Rotation Principle makes every move systemic, forcing constant holistic reasoning.
- **Duality of Phases:** The transition from an “architectural design” phase to a “tactical execution” phase tests the agent’s ability to make early decisions with irreversible consequences.

These characteristics, combined with the fact that it is a “contamination-free” environment (with no prior solutions in the LLMs’ training data), make it a powerful and precise tool for measuring the emergence of genuine reasoning.

## 7 Conclusion and Future Work

### 7.1 Conclusion

In this paper, we have dissected the reasoning process of an AI agent trained under the STAR-XAI Protocol. We have demonstrated that complex, strategic behavior can emerge from a simple, explicit set of rules when guided by a Socratic, interactive dialogue. Our central contribution is the identification of Second-Order Agency—the ability of an agent to reason about its own reasoning—as the key catalyst in this process.

Our analysis of the training logs shows a clear evolutionary path: from a “naive calculator” prone to “reasoning collapse” as described by Shojaee et al. (2025) [1], to a “disciplined strategist” capable of long-term planning and multi-objective optimization. We have argued that this transition is not a “black box” mystery but a transparent and replicable process, driven by the supervisor’s role in forcing the agent to deconstruct its own failures and co-create more robust cognitive protocols. Our findings challenge the notion that LRM reasoning is merely an “illusion”, suggesting instead that the limitations observed in other studies are often a consequence of non-agentic evaluation frameworks, an idea supported by critics like Khan et al. (2025) [3] and Song et al. (2025) [5]. The STAR-XAI protocol offers a practical methodology for overcoming this “agentic gap” and cultivating a reasoning process that is not only powerful but also auditable and trustworthy by design.

### 7.2 Future Work

The work presented here opens several avenues for future research. This paper is the second in a four-part series. The subsequent papers will provide a detailed examination of our “ante-hoc” transparency model in the context of XAI (Paper 3), and a breakdown of the “immune system” of protocols that ensures verifiable integrity and eliminates state-based hallucinations (Paper 4).

Beyond this series, the principles of inducing Second-Order Agency through dialogue are domain-agnostic. We plan to adapt this methodology to train and verify reliable agents in other high-stakes fields such as robotics and code generation.

Finally, the “Caps i Caps” environment itself, with its unique mechanics of indirect control and systemic consequences, has proven to be an exceptionally rich and “contamination-free” testbed for studying emergent reasoning. Future work will explore its potential as a tool for human cognitive training. This line of research is supported by a substantial body of evidence demonstrating that targeted training can enhance key cognitive abilities. For instance, meta-analyses have confirmed the malleability of spatial skills [6], and other work has established the advantages of working memory training [7]. Furthermore, the specific domain of strategy video games has been shown to enhance cognitive functions [9] and may even attenuate cognitive decline in older adults [8]. Investigating the anecdotal reports from expert players of a “cognitive expansion” sensation therefore becomes a promising research avenue, justifying the scaling of this work to the planned 8x8 and multiplayer configurations.

## References

- [1] Shojaee, P., Horton, M., Mirzadeh, I., Alizadeh, K., Bengio, S., & Farajtabar, M. (2025). The Illusion of Thinking: Understanding the Strengths and Limitations of Reasoning Models via the Lens of Problem Complexity. *arXiv preprint arXiv:2506.06941*.
- [2] Lawsen, A. (2025). The Illusion of the Illusion of Thinking: A Comment on Shojaee et al. (2025). *arXiv preprint arXiv:2506.09250*.
- [3] Khan, S., Madhavan, S., & Natarajan, K. (2025). A Comment On “The Illusion of Thinking”: Reframing the Reasoning Cliff as an Agentic Gap. *arXiv preprint arXiv:2506.18957*.
- [4] Dellibarda Varela, I., Romero-Sorozabal, P., Rocon, E., & Cebrian, M. (2025). Rethinking the Illusion of Thinking. *arXiv preprint arXiv:2507.01231*.
- [5] Song, Z., Yuet, S., & Zhang, J. (2025). Thinking Isn’t an Illusion: Overcoming the Limitations of Reasoning Models via Tool Augmentations. *arXiv preprint arXiv:2507.17699*.
- [6] Uttal, D. H., Meadow, N. G., Tipton, E., et al. (2013). The Malleability of Spatial Skills: A Meta-Analysis of Training Studies. *Psychological Bulletin*, 139(2), 352–402.
- [7] Moreau, D., & Conway, A. R. (2014). The case for a working memory training advantage. *Current Directions in Psychological Science*, 23(4), 277-282.
- [8] Basak, C., Boot, W. R., Voss, M. W., & Kramer, A. F. (2008). Can training in a real-time strategy video game attenuate cognitive decline in older adults? *Psychology and Aging*, 23(4), 765–777.
- [9] Oei, A. C., & Patterson, M. D. (2013). Enhancing cognition with video games: A multiple game training study. *PLoS ONE*, 8(3), e58546.