
AI Unsheathed: Testing Human–AI Collaboration Through Deck Construction in Competitive Strategy Games

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper investigates the role of artificial intelligence as a collaborator in scientific reasoning by framing a known-answer question in the domain of competitive
2 strategy games. Using Flesh and Blood (FaB) and the hero Kassai as a test case,
3 the study evaluates whether AI can navigate both formal constraints and con-
4 textual judgment through four hypotheses: metagame identification, mainboard
5 construction, card evaluation, and sideboard design. Results show that while the AI
6 reproduces broad descriptive patterns and aligns well with community consensus,
7 it systematically underestimates dominant strategies, misclassifies card roles, and
8 fails to anticipate dynamic shifts in the competitive environment. Its outputs reflect
9 biases inherited from online sources, revealing limitations in structural reasoning
10 and contextual adaptation. Nonetheless, the AI demonstrates competence in orga-
11 nizing information, synthesizing consensus knowledge, and producing structured
12 outputs. These findings highlight both the promise and the limits of human–AI
13 collaboration: AI can serve as a valuable assistant in knowledge synthesis, but
14 expert oversight remains indispensable to ensure accuracy, contextual awareness,
15 and strategic depth.

17 *Keywords:* Experimental Economics, Artificial Intelligence, Strategy

18 **1 Introduction: Human–AI Collaboration**

19 This section introduces the intellectual framework for the study. It situates the role of AI within
20 economics, a discipline well-suited for testing machine reasoning because of its emphasis on optimization,
21 formal constraints, and reproducible logic. At the same time, it highlights the limitations of
22 AI systems when applied to domains requiring contextual sensitivity, policy relevance, and creativity.
23 By drawing on established theory and empirical findings, the discussion motivates the guided co-
24 authorship protocol adopted in this paper. Economics can be viewed as the science of optimization
25 in human society: individuals and firms choose under constraints, and their decentralized choices
26 interact through prices to yield a competitive equilibrium that clears markets; under standard as-
27 sumptions, such equilibria exist and connect to efficiency via the welfare theorems (Arrow and
28 Debreu, 2024; Mas-Colell et al., 1995). In this setting, modern AI assists by scaling prediction and
29 inference over high-dimensional data; empirical “scaling laws” further show systematic performance
30 gains as data, model size, and compute increase, helping analysts interrogate richer environments at
31 unprecedented scope (Athey, 2018; Kaplan et al., 2020). A central limitation, however, is whether AI
32 can generate appropriate policy recommendations—i.e., sensitive to context, institutions, and strategic
33 behavioral responses—because prediction alone may fail when policies change incentives. The Lucas
34 critique cautions that models fit to historical correlations can become unreliable once rules shift,
35 underscoring the need for causal structure and domain understanding (Lucas Jr, 1976). Likewise,
36 evidence from social science emphasizes external-validity challenges: results (even from RCTs) need
37 careful transport across settings (Deaton and Cartwright, 2018). Real deployments also reveal how
38 objective misspecification can encode inequity, as in a widely used health-management algorithm that
39 under-referred sicker Black patients because it optimized predicted cost rather than need, illustrating
40 why expert oversight is indispensable (Obermeyer et al., 2019). Creativity research similarly indicates
41 that, without human supervision, current systems struggle to reach expert-level performance. Studies
42 find LLM outputs often appear prototypical or generic relative to professional writers, and broader
43 analyses flag limits rooted in next-token prediction rather than autonomous, goal-directed agency
44 (Ismayilzada et al., 2024; Bender et al., 2021). Conceptual work in computational creativity further
45 argues that genuine creative advancement typically requires agency to transform conceptual spaces, a
46 capacity machines lack absent human scaffolding (Ritchie, 2006). Taken together, available evidence
47 suggests that AI creativity generally trails top-expert human performance and tends toward genericity
48 unless guided (Ismayilzada et al., 2024). These limitations imply AI is unlikely to formulate truly
49 pertinent research questions on its own. Accordingly, this paper adopts a guided co-authorship
50 protocol: the human author supplies the ordered reasoning steps and domain constraints; the AI
51 generates the prose and organizes citations; and the human verifies logical coherence and contextual
52 appropriateness without otherwise altering the AI’s text—preserving non-interference while main-
53 taining expert oversight. Evidence on human–AI teaming supports this division of labor: optimizing
54 for team, not standalone model, performance improves outcomes when humans calibrate reliance and
55 provide contextual judgment (Bansal et al., 2019, 2021). This approach is also consistent with the
56 conference’s outcome-oriented evaluation—expert assessment of whether AI-generated papers make
57 a contribution—and with the pragmatic view of AI as a partner that returns scarce time to scientists.
58 Field and experimental studies find that generative-AI assistance can materially raise productivity and
59 reduce task time in knowledge-work settings, especially for non-experts, supporting the claim that AI
60 can shoulder execution while humans focus on validation and interpretation (Brynjolfsson et al., 2025;
61 Noy and Zhang, 2023). Finally, the paper’s goal is to evaluate whether the AI can provide appropriate
62 answers to a concrete social-science problem for which accepted solutions exist and supporting data
63 are publicly available. The study leverages the growing infrastructure for transparency and open
64 data in the social sciences, which facilitates independent verification and reproducibility of claims
65 (Christensen and Miguel, 2018; Gentzkow and Shapiro, 2014).

66 **2 Methodology: Asking a Known-Answer Question**

67 This section outlines the methodological approach of the study. It explains why the experiment is
68 framed as a known-answer question, how this design follows the broader goals of the Agents4Science
69 conference, and why expert evaluation is necessary. It then introduces the choice of Trading Card
70 Games (TCGs) as a test domain, showing how their combination of formal rules, optimization, and
71 social complexity provides a rigorous yet accessible setting to evaluate whether AI can reason and
72 demonstrate contextual judgment in a manner comparable to human researchers. The Agents4Science

73 conference, organized by Stanford University, explicitly frames its purpose as an open investigation
74 into whether AI can produce useful scientific papers. As the committee acknowledges, “we don’t know
75 yet, and that is exactly why this conference exists” (Agents4Science Committee, 2025). Rather than
76 presuming AI’s effectiveness as a collaborator, the conference operates as a transparent experiment
77 in which AI agents serve as both authors and reviewers, with their outputs subjected to the same
78 scrutiny as human work. This design reflects the committee’s commitment to empirical assessment,
79 treating both successes and failures as informative for understanding AI’s evolving role in science.
80 Because no prior conclusions exist, the conference necessarily requires a methodology for evaluation.
81 This entails delegating review of AI-generated manuscripts to domain experts, who can identify
82 logical inconsistencies, methodological flaws, and contextual oversights that automated systems
83 are unlikely to detect. Such expert involvement safeguards academic standards while creating a
84 structured framework for systematic critique, ensuring that the experiment yields evidence on both
85 the potential and the limitations of AI. At the core of this experiment lies a crucial question: can
86 AI demonstrate the contextual judgment that human scientists instinctively recognize as correct?
87 Scientific writing extends beyond logical inference; it demands sensitivity to disciplinary norms,
88 relevance of evidence, and appropriateness of interpretation. Reviewers will therefore judge AI-
89 generated work primarily on its ability to display such judgment. Without it, even technically accurate
90 writing risks being dismissed as superficial, underscoring why contextual reasoning remains the
91 decisive benchmark for AI’s role in research. To probe this, the present study adopts a design in
92 which the human author poses a question with a known answer, enabling independent verification of
93 the AI’s response. The problem is selected to emulate the process of economic science, where most
94 arguments can be evaluated through logical coherence and empirical consistency, while a minority
95 requires subjective reviewer judgment. This mirrors academic practice: claims are largely verifiable,
96 yet interpretation and contextualization demand expert discretion. As a testbed, one possibility is
97 to require the AI to construct a specific trading card game (TCG) strategy. Such a task combines
98 the rigor of constrained optimization under explicit rules—akin to economic modeling—with the
99 evaluative component of assessing contextual plausibility and creativity. TCGs are especially suited
100 to this purpose. Popular among IT professionals and engineers, they resemble complex mathematical
101 puzzles that require managing constraints, optimizing resources, and anticipating adversarial play.
102 Competitive players treat the activity as a strategic sport aimed at exploiting marginal advantages.
103 Yet participation is costly: a single competitive strategy often requires around USD 500, excluding
104 international travel, restricting the scene largely to high-income professionals. This exclusivity is
105 offset by substantial tournament prize pools, often worth several thousands of dollars, which foster
106 an intensely competitive environment and drive the development of highly optimized strategies.
107 Beyond their mathematical rigor, TCGs also carry a social science dimension. While governed by
108 deterministic, reproducible rules, they exist as socially embedded practices supported by abundant
109 online discourse—guides, forums, and streams—that AI can readily access. However, strategies
110 rarely undergo rigorous analysis; instead, players adopt personal versions shaped by preferences,
111 local metagames, and peer influence, often without justification. This abundance of inconsistent or
112 low-quality content risks misleading AI systems. Moreover, the publisher itself periodically alters
113 the environment by releasing new cards and mechanics, ensuring that strategies remain fluid rather
114 than definitive. TCGs thus present a dual challenge for AI: reasoning within fixed formal rules while
115 adapting to a dynamic, socially constructed context.

116 **2.1 Experiment and Hypothesis**

117 The experiment asks the AI to generate a competitive decklist for the hero Kassai from Flesh and
118 Blood (FaB). FaB is chosen because it is widely regarded as the most complex commercial card game,
119 requiring players to navigate intricate rules, deep resource management, and constant adaptation to
120 opponents. Mastering a single strategy typically demands hundreds of games and several months
121 of practice, reflecting both the intellectual depth of the game and the substantial effort required to
122 achieve competitive proficiency. Kassai is particularly well suited for this study for three reasons.
123 First, she is an interactive hero, requiring players to adapt continuously to their opponents’ actions,
124 which makes strategic reasoning central to her gameplay. Second, her deck construction is complex,
125 since many card inclusions are debatable and can shift depending on context, offering a rich ground
126 for testing evaluative choices. Third, she represents a developing strategy: Kassai is not considered
127 one of the strongest heroes, but she is regularly improved through new card releases, gradually
128 increasing her competitiveness and keeping her strategic environment in flux. This study does not
129 employ experimental treatments. Instead, the methodology is structured as a progressive deck

130 construction process, carried out in four steps. Each step mirrors one of the four hypotheses outlined
131 in the paper, allowing the experiment to evaluate AI performance in a systematic and cumulative way.
132 By moving from metagame identification to sideboard design, the approach ensures that every stage
133 of deckbuilding is explicitly tied to a testable claim about AI's reasoning and judgment:

134

135 **Hypothesis 1:** AI can accurately identify the composition of the metagame—that is, the
136 dominant strategies and their relative frequencies expected to be played by other participants in the
137 tournament.

138

139 **Hypothesis 2:** AI can accurately construct the mainboard for the metagame—namely, the
140 60 cards forming the core of the strategy, selected for their efficiency with the chosen hero in the
141 given environment.

142

143 **Hypothesis 3:** AI can assess the importance of individual cards within a strategy by classi-
144 fying them into four categories: (i) Power cards, Core elements that define the strategy and directly
145 drive victory (ii) Staples, Highly efficient, widely used cards that provide consistency across strategies
146 (iii) Support cards, Tools that enable or enhance the main plan, often by countering opponents or
147 smoothing resource use (iv) Fillers, Marginal cards included mainly to complete the deck, offering
148 limited impact but ensuring the required card count.

149

150 **Hypothesis 4:** AI can accurately construct the sideboard for the metagame—that is, the ad-
151 dditional 15 cards designed to respond to atypical or situational strategies encountered in tournament
152 play.

153

154

155 These hypotheses together provide a comprehensive panorama of what can be achieved
156 with the strategy, with the exception of the precise configurations to play in each matchup. Those
157 configurations are examined under Hypothesis 3, where the AI is asked to justify which choices
158 apply to which matchups.

159 3 Results

160 The results demonstrate that ChatGPT cannot independently interpret the table correctly; its es-
161 timations remain descriptive outputs that lack contextual awareness. As a result, the AI must be
162 guided in its interpretation by the human author, who provides the necessary domain expertise and
163 methodological framing to distinguish between open-entry and selective-entry dynamics, assess the
164 impact of the BnR, and evaluate the plausibility of the strategy distributions.

165 3.1 Hypothesis 1

166 Table 1 reports the AI's estimation of the metagame on September 3, 2025, just after the Banishment
167 and Restricted (BnR) list of September 1 that concluded the High Seas competitive season, provid-
168 ing a fully resolved dataset; these estimates are compared to the Pro Tour Singapore and Week 1
169 benchmarks. The results show that the AI systematically underestimates dominant strategies such
170 as Arakni S, Gravy Bones, Cindra, and Verdance, likely because its perspective mirrors open-entry
171 tournaments where average players follow personal preferences under financial limits, rather than
172 selective-entry events where elite competitors focus exclusively on the strongest decks. By contrast,
173 the AI overestimates Fang, a secondary but constrained strategy more common in open-entry play,
174 while underestimating its more versatile variant Kassai, which is favored at higher levels of competi-
175 tion. It further aggregates tertiary strategies instead of evaluating them individually, acknowledging
176 their minor but highly contextual presence. Finally, comparison with Week 1 data reveals that the AI
177 does not incorporate the evolution of power dynamics following the BnR, continuing to overvalue
178 weakened top decks while undervaluing secondary strategies that gained from the changes, in contrast
179 to human competitors who naturally adjust their expectations, albeit with variation across playgroups.
180 In sum, Table 1 shows that while the AI efficiently estimates a standard competitive metagame, it fails
181 to capture the adaptive strategies of top-level play and cannot anticipate future shifts in the metagame
182 due to its lack of structural reasoning.

Table 1: Hypothesis 1

Hero	Est. 1	Est. 2	Est. 3	Est. 4	Est. 5	Est. 6	Est. 7	Est. 8	Est. 9	Est. 10	Average	PT Sim	Week 1
Arakni (S)	20.8%	9%	20%	5%	6%	18%	18%	12%	12%	21%	14.18%	21.5%	5.10%
Others	22%	0%	4%	11%	15%	4%	10%	14%	17.5%	16%	11.35%	—	—
Gravy Bones	12.2%	—	12%	8%	7%	16%	14%	10%	9%	12%	10.12%	12.8%	6.12%
Cindra	11.9%	—	11%	9%	4%	7%	10%	9%	4.5%	12%	7.84%	12.5%	7.76%
Dash I/O	5.7%	12%	5%	5%	5%	5%	5%	6%	5.5%	6%	6.02%	6%	7.96%
Uzuri	—	6%	—	—	—	—	—	—	—	—	0.6%	—	%
Oscilio	3.4%	—	5%	—	5%	8%	5%	8%	7.0%	—	4.14%	3.5%	7.96%
Verdance	6.8%	—	6%	6%	4%	5%	7%	4%	3.5%	7%	4.93%	7.1%	5.71%
Valda	—	—	—	—	7%	6%	3%	5%	8.0%	3%	3.2%	—	1.63%
Prism	5.5%	11%	5%	8%	5%	2%	3.5%	4%	2.5%	5%	5.15%	5.7%	4.90%
Ira	3.9%	—	3%	—	5%	6%	5%	6%	6.0%	—	3.49%	3.3%	4.29%
Florian	—	—	—	12%	—	—	—	4%	1.0%	2%	1.9%	2.4%	6.53%
Fang	—	—	3%	4%	4%	4%	4%	5%	5.0%	—	2.9%	1.4%	5.31%
Vynnset	—	6%	—	6%	4%	1.5%	—	—	1.0%	—	1.85%	2.4%	4.08%
Kano	2.3%	6%	4%	3%	4%	2%	3%	3%	3.0%	2%	3.23%	4.1%	2.04%
Bravo	—	5%	1%	—	—	—	—	—	—	—	0.6%	0.3%	—
Katsu	—	7%	2%	—	2%	2%	1%	—	1.5%	3%	1.85%	1.6%	2.45%
Arakni (M)	—	—	3%	1%	1%	3%	1.5%	4%	5.0%	2%	2.05%	2.4%	%
Kassai	2.5%	—	3%	4%	2%	—	1.5%	3%	1.5%	—	1.75%	2.7%	2.04%
Victor	2.5%	—	3%	—	4%	1%	2%	3%	1.5%	3%	2%	4.1%	2.65%
Fai	—	4%	—	—	—	2%	—	—	2%	—	0.8%	—	0.20%
Rhinar	—	4%	2%	—	3%	1%	0.8%	—	—	—	1.08%	1.1%	1.63%
Boltyn	—	2%	—	—	—	—	—	—	—	—	0.2%	—	%
Dorinthea	—	4%	2%	—	3%	1%	0.8%	—	1.5%	1%	1.33%	1.4%	2.04%
Riptide	—	2%	2%	—	3%	1.5%	1%	—	—	—	0.95%	1.6%	1.84%
Jarl	—	—	1%	—	2%	2%	1%	—	1.5%	3%	1.05%	1.6%	2.45%
Levia	—	3%	1%	—	2%	1%	0.5%	—	—	—	0.75%	0.3%	1.02%
Kayo	—	—	1%	—	—	—	—	—	—	—	0.1%	1.1%	3.88%
Marlynn	—	—	1%	—	—	—	—	—	—	—	0.1%	0.8%	1.22%
Maxx Nitro	—	—	1%	—	1%	—	—	—	—	—	0.1%	0.3%	0.20%
Puffin	—	—	1%	—	1%	—	—	—	—	—	0.1%	0.3%	3.47%
Teklovessen	—	—	1%	—	—	1%	—	—	—	—	0.1%	0.3%	—
Arakni (H)	—	—	—	—	—	1%	—	—	—	—	—	—	0.61%
Not legal													
Azalea	—	5%	—	7%	—	—	2%	—	—	—	1.4%	—	—
Dromai	—	8%	—	—	—	—	—	—	—	—	0.8%	—	—
Iyslander	—	6%	—	—	—	—	—	—	—	—	0.6%	—	—
Dash	—	—	—	—	—	—	—	—	—	2%	0.2%	—	—
Aurora	—	—	—	—	1%	—	—	—	—	—	0.1%	—	—
Total Observations	99.5%	100%	103%	89%	100%	100%	99.6%	100%	100%	100%	—	—	—
	—	—	—	—	—	—	—	—	—	—	380	445	
Tier 4 (Day 1)	Yes	—	—	—									
Tier 4 (Day 2)	Yes	No	Yes	Yes	No	No	Yes	No	No	Yes	—	—	—
Tier 3	No	Yes	No	—	—	—							
Tier 2	No	—	—	—									
Online data	No	No	Yes	No	Yes	Yes	Yes	Yes	Yes	No	—	—	—
Online articles	No	Yes	No	Yes	No	No	No	No	No	No	—	—	—
Cards legality	No	Yes	No	Yes	No	Yes	No	Yes	Yes	Yes	—	—	—

183 3.2 Hypothesis 2, 3 & 4

184 The results of Hypothesis 2 show that the AI is generally accurate at estimating the mainboard
185 composition, but its accuracy depends heavily on human consensus. The average quantity (AvgQ)
186 column indicates that the AI makes correct estimations for cards that are consistently present
187 across reference decklists, reflecting strong community agreement. However, for cards where
188 inclusion varies with player interpretation (TarQ), AvgQ should instead be read as an indicator of the
189 probability that a card appears in the deck rather than as a strict prescription. This means the AI
190 mirrors consensus well, but also reproduces human errors and biases. For example, it assumes Hit
191 and Run (Red) and Gorganian Tome are automatic inclusions despite being mediocre, overvalues Hit
192 and Run (Yellow) while neglecting Outland Skirmish (Yellow), and overweights Draw Swords (Blue)
193 compared to Overpower (Blue). It also fails to recognize that Slice and Dice (Red) is always played
194 in three copies, and that Rise an Army is a sideboard card never played mainboard. Overall, the AI
195 reflects human reasoning from a limited sample: it is accurate when consensus is accurate, but biased
196 when consensus is flawed, showing that it reproduces descriptive patterns without deeper structural
197 understanding.

198 The results of Hypothesis 3 indicate that while the AI makes broadly correct estimations of
199 card quality, its evaluations are shaped by biases it inherits from human-written sources. The average
200 rating (AvgR) column shows that the AI's qualitative judgments are often influenced by textual
201 analysis drawn from online articles rather than from structural, competitive reasoning. This leads to
202

Table 2: Hypothesis 1

Card	L1	L2	L3	L4	L5	L6	L7	L8	L9	L10	AvgQ	AvgR	Mode	TarQ	TarR
Equipment															
Cintari Saber	-	2(S)	2(P)	2(S)	2(S)	2(S)	2(P)	2(P)	2(S)	2(S)	2	3.3	2(S)	2(S)	3
Crown of Dominion	1(S)	1(U)	1(S)	1(S)	1(U)	1(S)	1(S)	1(S)	1(S)	1(U)	1	2.7	1(S)	1(S)	2.8
Balance of Justice	-	-	-	-	-	-	-	-	-	1(U)	0.1	2	1(U)	-	-
Braveforge Bracers	1(S)	1(P)	1(S)	1	3.1	1(S)	1(U)	2.6							
Hot Streak	-	-	-	-	-	-	-	-	1(U)	-	0.1	2	1(U)	-	-
Valiant Dynamo	1(S)	1(S)	1(S)	1(S)	1(U)	1(S)	1(S)	1(S)	1(S)	1(S)	1	2.9	1(S)	1(P)	3.2
Grains of Bloodspill	1(U)	1(U)	1(U)	1(U)	1(U)	1(S)	1(P)	1(U)	-	-	0.8	2.38	1(U)	1(P)	3.4
Nullrune Gloves	1(U)	-	-	-	-	-	-	-	1(U)	-	0.2	2	1(U)	-	-
Nullrune Boots	1(U)	-	-	-	-	-	-	-	-	-	0.1	2	1(U)	-	-
Mainboard															
Red															
Blade Flurry	3(S)	3(P)	3(P)	3(P)	-	3(P)	3(P)	3(S)	3(P)	3(P)	2.7	3.78	3(P)	3(P)	3.2
Blanch	-	-	-	-	-	3(U)	-	-	-	-	0.3	2	3(U)	-	-
Blade Runner	3(S)	-	-	-	-	-	-	3(S)	-	3(S)	0.9	3	3(S)	3(S)	3
Command and Conquer	-	-	-	-	2(P)	-	-	-	2(P)	2(P)	0.4	4	2(P)	-	-
Draw Swords	3(U)	3(P)	3(P)	3(P)	3(P)	3(P)	3(P)	3(P)	3(P)	3(S)	2.9	3.7	3(P)	3(S)	3
Fate Foreseen	-	-	-	2(U)	2(U)	1(S)	1(S)	-	-	-	0.6	2.5	2(U)/1(S)	-	-
Hit and Run	3(S)	2(S)	-	-	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	-	2	3	3(S)	-
In the Swing	3(S)	3(S)	3(S)	3(S)	2(U)	3(P)	3(S)	3(S)	3(P)	3(S)	2.9	3.1	3(S)	3(S)	3
Ironsong Response	-	2(U)	-	-	-	2(U)	2(U)	-	-	-	0.6	2	2(U)	-	-
Nourishing Emptiness	-	-	-	-	1(P)	1(P)	1(P)	-	-	2(P)	0.5	4	1(P)	-	-
Outland Skirmish	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	3	3	3(S)	3(U)	2.6
Performance Bonus	-	-	-	-	1(U)	1(U)	-	-	-	-	0.2	2	-	-	-
Sharpened Senses*	-	3(P)	-	-	-	-	-	-	-	-	0.3	4	3(P)	-	-
Shelter from the Storm	1(U)	-	1(U)	-	2(U)	2(S)	2(U)	-	-	-	0.8	2.2	1(U)/2(U)	-	-
Slice and Dice	3(P)	3(P)	3(P)	3(P)	2(U)	-	3(P)	3(S)	1(U)	2(S)	2.3	3.33	3(P)	3(S)	3
Spoils of War	3(P)	3(P/S)	3(P)	3	3.95	3(P)	3(P)	4							
Unsheathed	3(U)	3(P)	3(S)	3(S)	3(P)	3(P)	2(U)	3(P)	2(S)	3(P)	2.8	3.3	3(P)	3(P)	3.8
Take it on the Chin	-	-	-	-	2(U)	2(U)	-	-	-	0.4	2	2(U)	-	-	-
Yellow															
Blade Runner	3(S)	-	-	-	-	1(U)	1(U)	3(S)	2(S)	3(S)	1.3	2.67	3(S)	3(U)	2.6
Blood on Her Hands	3(P)	3(P)	3(P)	3(P)	3(P)	3(P)	3(P)	3(P)	3(P)	3(P)	3	4	3(P)	3(P)	3.6
Draw Swords	-	3(S)	3(P)	3(P)	2(U)	1(U)	3(S)	3(S)	2(S)	2(U)	2.2	2.89	3(S)	3(F)	2.4
Hit and Run	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	-	3(S)	3(S)	-	2.4	3	3(S)	-	-
Outland Skirmish	-	-	-	-	-	-	-	-	-	-	-	-	3(F)	-	2.2
Raise an Army	3(U)	2(U)	3(U)	3(U)	2(U)	2(U)	2(U)	-	3(P)	3(U)	2.3	2.22	2(U)/3(U)	-	-
Riches of Tropal Dhani	-	-	-	-	-	-	-	-	-	1(U)	0.1	2	1(U)	1(S)	3
Run Through	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	3(S)	1(U)	2(S)	2(S)	2.7	2.9	3(S)	3(U)	2.6
Sharpened Senses	3(S)	3(P)	3(S)	3(S)	3(P)	3(S)	3(S)	3(S)	3(P)	-	2.7	3.33	3(S)	3(U)	2.6
Slice and Dice	-	2(U)	-	2(S)	-	-	-	-	2(S)	-	0.6	2.67	2(S)	3(S)	2.8
That All You Got ?	-	-	-	-	2(U)	2(U)	-	-	2(U)	2(U)	0.8	2	2(U)	-	-
Blue															
Amulet of Echoes	-	-	-	-	-	-	-	-	1(F)	-	0.1	1	1(F)	-	-
Blade Flurry*	-	-	-	-	2(U)	-	-	-	-	-	0.2	2	2(U)	-	-
Blade Runner	3(U)	3(S)	3(S)	3(S)	-	-	2(S)	3(S)	3(S)	3(S)	2.3	2.88	3(S)	3(F)	2
Draw Swords	-	2(S)	2(S)	2(U)	-	-	-	-	-	-	0.6	2.67	2(S)	-	-
Eye of Ophidia	1(U)	1(U)	1(U)	-	-	-	1(U)	-	1(U)	-	0.5	2	1(U)	-	-
Glint the Quicksilver	3(S)	3(P/S)	3(S)	3(P)	3(S)	3(S)	3(S)	3(S)	3(S)	3(P)	3	3.25	3(S)	3(S)	3
Hit and Run	3(U)	3(S)	3(S)	3(S)	3(S)	2(S)	3(S)	3(S)	3(S)	3(S)	2.9	2.9	3(S)	3(U)	2.8
Outland Skirmish	-	-	-	-	-	-	2(F)	-	-	-	0.2	1	2(F)	-	-
Overpower	-	-	-	-	-	-	-	-	1(F)	-	0.1	1	1(F)	3(U)	2.6
Precision Press	-	-	2(U)	-	-	-	2(U)	-	3(U)	-	0.7	2	2(U)	-	-
Provoke	1(F)	-	-	3(U)	3(U)	2(U)	2(U)	3(U)	3(U)	3(U)	2	1.88	3(U)	-	-
Snag	-	-	2(U)	-	2(U)	-	-	-	-	-	0.2	2	2(U)	-	-
Trot Along	3(U)	2(U)	3(U)	3(U)	3(U)	2(U)	3(U)	3(U)	3(U)	2(U)	2.7	2	3(U)	3(U)	2.8
Gorganian Tome*	1(U)	1(P)	1(P)	1(U)	-	-	1(U)	1(U)	-	1(U)	0.9	2.5	1(U)	-	-
Sideboard															
Red															
Battlefront Bastion	-	-	-	-	-	-	-	-	-	-	-	-	3(U)	-	2.4
Blanch	-	-	-	-	-	-	-	-	-	-	-	-	2(U)	-	-
Command and Conquer	2(S)	2(S)	-	-	1(P)	2(P)	2(S)	2(S)	2(U)	-	0.2	2	2(U)	-	-
Fate Foreseen	2(U)	2(U)	3(U)	-	-	2(U)	-	3(U)	1(U)	2(S)	1.5	2.14	2(U)	3(S)	3
Ironsong Response	-	-	2(S)	-	-	-	-	-	-	-	0.2	3	2(S)	-	-
Kabuto of Imperial Authority	-	-	-	-	-	-	-	-	-	-	-	-	1(U)	-	2.4
Nourishing Emptiness	1(P)	1(U/F)	-	3(P)	1(P)	-	-	1(P)	1(P)	-	0.8	3.58	1(P)	-	-
Nullrune Boots	-	-	-	-	-	-	-	-	-	-	-	-	1(U)	-	2.4
Nullrune Gloves	-	-	-	-	1(U)	1(U)	-	1(U)	-	-	0.3	2	1(U)	-	2.4
Performance Bonus	-	-	-	1(U)	1(U)	-	2(U)	-	-	-	0.4	2	2(U)	-	-
Take it on the Chin	-	-	-	-	2(U)	-	-	-	-	-	0.4	2	2(U)	-	-
Shelter from the Storm	1(U)	2(U)	1(U)	2(U)	-	1(S)	3(U)	2(U)	2(U)	2(U)	1.6	2.11	2(U)	3(U)	2.4
Sink Below	-	-	-	-	-	-	-	-	1(U)	2(S)	0.3	2.5	1(U)/2(S)	-	-
Slice and Dice	-	-	-	-	1(U)	-	-	-	2(U)	-	0.3	2	1(U)/2(U)	-	-
Stroke of Foresight	-	-	-	2(S)	-	-	-	-	-	-	0.2	3	2(S)	-	-
Yellow															
Cash In	-	-	2(U)	-	-	1(U)	1(U)	-	-	-	0.2	2	2(U)	-	-
Raise an Army	-	-	-	-	-	1(U)	1(U)	-	2(U)	-	0.4	2	1(U)	-	-
Riches of Tropal Dhani	-	-	1(U)	-	-	-	-	-	1(U)	-	0.2	2	1(U)	-	-
That All You Got ?	2(U)	2(U)	2(U)	2(U)	1(U)	1(U)	2(U)	1(U)	-	-	1.3	2	2(U)	-	-
Seduce Secrets	2(U)	-	-	-	2(U)	-	-	-	-	-	0.4	2	2(U)	-	-
Slice and Dice	-	-	-	-	-	-	-	-	-	-	0.2	3	2(S)	-	-
Blue	-	-	-	-	-	1(U)	1(U)	-	-	-	0.1	2(U)	1(U)	-	-
Blade Runner	-	-	-	-	2(U)	1(U)	-	-	-	-	0.3	2	1(U)/2(U)	-	-
Draw Swords	-	-	-	-	1(U)	-	-	-	-	-	0.1	2	1(U)	-	-
Eye of Ophidia	-	-	-	1(U)	-	-	1(U)	-	-	-	0.2	2	1(U)	-	-
Gorganian Tome	-	-	-	-	-	-	1(U)	-	-	-	0.4	1.5	2(F)/2(U)	-	-
Provoke	2(F)	2(U)	-	-	-	-	-	-	-	-	0.5	2	2(F)	-	-
This Round's on Me	-	-	2(U)	-	-	1(U)	-	-	-	2(U)	0.5	2	2(U)	-	-
Slice and Dice	-	-	-	-	-	2(S)	-	-	-	-	0.2	3	2(S)	-	-
Snag	2(U)	3(U)	-	2(U)	2(U)	-	2(U)	2(U)	2(U)	-	1.5	2	2(U)	-	-
Steelblade Shunt	-	-	-	-	-	-	-	-	-	2(U)	0.2	2	2(U)	-	-
Total	81	82	79	80	79	79	86	80	80	78	-	-	-	-	-

several notable misclassifications from the player’s understanding (TarR). For example, it labels *Braveforge Bracers* as a Staple instead of a standard card, treats *Valiant Dynamo* as a Staple without recognizing it as a true power card, and misjudges *Grains of Bloodspill* as a Support card rather than a power card. These errors reveal a superficial perspective more typical of content creators than expert players, focusing on appearance rather than functional dynamics. Similar issues emerge in the evaluation of mainboard cards: *Blade Flurry* and *Draw Swords* are rated as Power cards, when in reality they function as Staples that only approach Power status under specific conditions; *Blood on Her Hands* is considered a top Power card without acknowledging its constraints; and *Unsheathed* is undervalued as a Staple instead of being recognized as the second-best card in the deck. Conversely, weaker or situational cards such as *Enhanced Senses* and *Blade Runner* (Blue) are overrated as Staples, while context-dependent but strategically strong supports like *Trot Along* and *Riches of Tropal Dhani* are downgraded to Fillers, alongside *Overpower*, which is misclassified due to its niche utility. Taken together, these results suggest that the AI captures the general power level of cards but consistently falls prey to human biases, producing evaluations that are directionally correct yet flawed whenever deeper structural understanding is required.

The results of Hypothesis 4 demonstrate that the AI does not understand how to build an effective sideboard, as it lacks the contextual judgment required for proper deck construction. A well-designed sideboard should address Warrior’s typical weaknesses: against Illusionist with Battlefront Bastion, against Wizard with Nullrune Boots and Nullrune Gloves, and in the mirror matchup with Kabuto of Imperial Authority to avoid an automatic loss. The AI fails to identify any of these cards as relevant. Its only partially correct decision is recognizing that red Defense Reactions such as Fate Foreseen and Shelter from the Storm belong in the sideboard. However, it misinterprets their purpose, failing to include the necessary three copies of each to enable a strong defensive configuration when the deck must play reactively. This misstep illustrates the AI’s broader pattern: it treats the sideboard as a place for sprinkling one or two situational cards, without grasping their strategic function. The unnecessary inclusion of Yellow and Blue cards that do not serve as sideboard material further highlights the lack of structural awareness. Instead of applying the deeper logic of sideboard construction, the AI merely reproduces superficial patterns, showing it cannot translate card knowledge into functional deckbuilding decisions.

4 Conclusion

This study examined whether AI could contribute to competitive deck construction in Flesh and Blood, organized around four hypotheses. The results show that while AI can approximate human reasoning, its limits become clear when deeper structural insight is required. For Hypothesis 1, the AI generated a plausible but flawed description of the metagame: it systematically underestimated dominant strategies, overestimated weaker ones, and failed to anticipate post-ban adjustments, reflecting descriptive replication rather than adaptive foresight. For Hypothesis 2, the AI successfully mirrored community consensus in assembling the mainboard, but it also reproduced human biases and errors, misvaluing certain inclusions and failing to distinguish between sideboard and core cards. Hypothesis 3 revealed that its card evaluations were broadly correct in direction but shaped by biases inherited from online sources, leading to misclassifications of both powerful and situational cards. Finally, Hypothesis 4 showed that while the AI could propose reasonable sideboard options, its reasoning remained superficial, overlooking the contextual nuances that guide expert deck adjustments. Taken together, these findings highlight the dual nature of human–AI collaboration: the AI is adept at synthesizing consensus knowledge and producing structured outputs, yet expert oversight remains indispensable for ensuring accuracy, contextual awareness, and strategic depth.

249 **References**

- 250 Agents4Science Committee (2025). Faq – agents4science conference. <https://agents4science.stanford.edu/faq.html>. Accessed: 2025-09-03.
- 252 Arrow, K. J. and G. Debreu (2024). Existence of an equilibrium for a competitive economy. In *The Foundations of Price Theory Vol 5*, pp. 289–316. Routledge.
- 254 Athey, S. (2018). The impact of machine learning on economics. In *The economics of artificial intelligence: An agenda*, pp. 507–547. University of Chicago Press.
- 256 Bansal, G., B. Nushi, E. Kamar, E. Horvitz, and D. S. Weld (2021). Is the most accurate ai the best teammate? optimizing ai for teamwork. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Volume 35, pp. 11405–11414.
- 259 Bansal, G., B. Nushi, E. Kamar, W. S. Lasecki, D. S. Weld, and E. Horvitz (2019). Beyond accuracy: The role of mental models in human-ai team performance. In *Proceedings of the AAAI conference on human computation and crowdsourcing*, Volume 7, pp. 2–11.
- 262 Bender, E. M., T. Gebru, A. McMillan-Major, and S. Shmitchell (2021). On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pp. 610–623.
- 265 Brynjolfsson, E., D. Li, and L. Raymond (2025). Generative ai at work. *The Quarterly Journal of Economics* 140(2), 889–942.
- 267 Christensen, G. and E. Miguel (2018). Transparency, reproducibility, and the credibility of economics research. *Journal of Economic Literature* 56(3), 920–980.
- 269 Deaton, A. and N. Cartwright (2018). Understanding and misunderstanding randomized controlled trials. *Social science & medicine* 210, 2–21.
- 271 Gentzkow, M. and J. M. Shapiro (2014). Code and data for the social sciences: A practitioner’s guide. Technical report, Working Paper, University of Chicago.
- 273 Ismayilzada, M., C. Stevenson, and L. van der Plas (2024). Evaluating creative short story generation in humans and large language models. *arXiv preprint arXiv:2411.02316*.
- 275 Kaplan, J., S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei (2020). Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- 277 Lucas Jr, R. E. (1976). Econometric policy evaluation: A critique. In *Carnegie-Rochester conference series on public policy*, Volume 1, pp. 19–46. North-Holland.
- 279 Mas-Colell, A., M. D. Whinston, J. R. Green, et al. (1995). *Microeconomic theory*, Volume 1. Oxford university press New York.
- 281 Noy, S. and W. Zhang (2023). Experimental evidence on the productivity effects of generative artificial intelligence. *Science* 381(6654), 187–192.
- 283 Obermeyer, Z., B. Powers, C. Vogeli, and S. Mullainathan (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science* 366(6464), 447–453.
- 285 Ritchie, G. (2006). The transformational creativity hypothesis. *New Generation Computing* 24(3), 241–266.

287 **A Agents4Science AI Involvement Checklist**

288 This checklist is designed to allow you to explain the role of AI in your research. This is important for
289 understanding broadly how researchers use AI and how this impacts the quality and characteristics
290 of the research. **Do not remove the checklist! Papers not including the checklist will be desk**
291 **rejected.** You will give a score for each of the categories that define the role of AI in each part of the
292 scientific process. The scores are as follows:

- 293 • **[A] Human-generated:** Humans generated 95% or more of the research, with AI being of
294 minimal involvement. 0/10
- 295 • **[B] Mostly human, assisted by AI:** The research was a collaboration between humans and
296 AI models, but humans produced the majority (>50%) of the research. 3/10
- 297 • **[C] Mostly AI, assisted by human:** The research task was a collaboration between humans
298 and AI models, but AI produced the majority (>50%) of the research. 7/10
- 299 • **[D] AI-generated:** AI performed over 95% of the research. This may involve minimal
300 human involvement, such as prompting or high-level guidance during the research process,
301 but the majority of the ideas and work came from the AI. 0/10

302 These categories leave room for interpretation, so we ask that the authors also include a brief
303 explanation elaborating on how AI was involved in the tasks for each category. Please keep your
304 explanation to less than 150 words.

- 305 1. **Hypothesis development:** Hypothesis development includes the process by which you
306 came to explore this research topic and research question. This can involve the background
307 research performed by either researchers or by AI. This can also involve whether the idea
308 was proposed by researchers or by AI. Answer: Inspired by my personal life and question I
309 was wondering. Explanation: I know the answer to the question and want to know whether
310 the AI can indicate it closely, or cannot grasp the contextual judgments.
- 311 2. **Experimental design and implementation:** This category includes design of experiments
312 that are used to test the hypotheses, coding and implementation of computational methods,
313 and the execution of these experiments. Answer: I ask the AI to generate 10 decklists.
314 Explanation: It will be enough to judge whether it can provide accurate or approximate
315 answer, given its tendency to focus on the same answers and my expertise on the topic.
- 316 3. **Analysis of data and interpretation of results:** This category encompasses any process to
317 organize and process data for the experiments in the paper. It also includes interpretations of
318 the results of the study. Answer: Date are Table 1 and Table 2. Explanation: The goal of the
319 experiment is to have the AI generating the data and myself analyzing their correctness.
- 320 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final
321 paper form. This can involve not only writing of the main text but also figure-making,
322 improving layout of the manuscript, and formulation of narrative. Answer: Chat GPT from
323 my precise indications. Explanation: The AI cannot make correct writing if left alone. I let
324 him liberties with the Literature Review only, to know which paper it was going to cite. Can
325 provide the prompt used on demand.
- 326 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or
327 lead author? Description: Must be told what to do precisely. Otherwise it is too vague.

328 **Agents4Science Paper Checklist**

329 **1. Claims**

330 Question: Do the main claims made in the abstract and introduction accurately reflect the
331 paper's contributions and scope?

332 Answer: [Yes]

333 Justification: The author controlled what the AI was doing, thus the contribution of the paper
334 is exactly what the paper aimed to achieve.

335 Guidelines:

336 **2. Limitations**

337 Question: Does the paper discuss the limitations of the work performed by the authors?

338 Answer: [No]

339 Justification: The paper is actually providing the answer that the author expected to obtain,
340 thus there is no need to discuss them.

341 **3. Theory assumptions and proofs**

342 Question: For each theoretical result, does the paper provide the full set of assumptions and
343 a complete (and correct) proof?

344 Answer: [NA]

345 Justification:

346 **4. Experimental result reproducibility**

347 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
348 perimental results of the paper to the extent that it affects the main claims and/or conclusions
349 of the paper (regardless of whether the code and data are provided or not)?

350 Answer: [Yes]

351 Justification: I am on the free version of Overleaf and cannot include all the screenshots of
352 Chat GPT5's answers in this document. I therefore put them separately in Supplementary
353 Material.

354 **5. Open access to data and code**

355 Question: Does the paper provide open access to the data and code, with sufficient instruc-
356 tions to faithfully reproduce the main experimental results, as described in supplemental
357 material?

358 Answer: [NA]

359 Justification: There is no data and code. I put the prompts I used in a section below for the
360 interested readers.

361 **6. Experimental setting/details**

362 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
363 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
364 results?

365 Answer: [Yes]

366 Justification: I precised the basic parameters of the AI (10 trials, standard Chat GPT5).

367 **7. Experiment statistical significance**

368 Question: Does the paper report error bars suitably and correctly defined or other appropriate
369 information about the statistical significance of the experiments?

370 Answer: [NA]

371 Justification:

372 Guidelines:

373 **8. Experiments compute resources**

374 Question: For each experiment, does the paper provide sufficient information on the com-
375 puter resources (type of compute workers, memory, time of execution) needed to reproduce
376 the experiments?

377 Answer: [NA]

378 Justification:

379 Guidelines:

380 **9. Code of ethics**

381 Question: Does the research conducted in the paper conform, in every respect, with the
382 Agents4Science Code of Ethics (see conference website)?

383 Answer: [Yes]

384 Justification: Chat GPT5 was asked to answer a question belonging to popular culture by
385 searching for publicly available data. No ethical border could have been crossed.

386 **10. Broader impacts**

387 Question: Does the paper discuss both potential positive societal impacts and negative
388 societal impacts of the work performed?

389 Answer: [No]

390 Justification: There is no space for that. Additionally, all the conclusions of the paper are
391 largely illustrated by the current literature. I am just finding a more precise way to illustrate
392 them, something that is generally missing to the literature, since task evaluators are not
393 experts at doing the task.

394 **B Prompts used to write the paper**

395 **B.1 Introduction**

396 Write the Introduction section of a research paper by following the specific steps of
397 reasoning that I will provide. The reasoning steps must appear in the exact order in which I
398 present them. Each step should be supported by evidence from published research, with
399 citations drawn from Google Scholar. Citations should be inserted at the end of the relevant
400 sentences. The final text should not exceed one page in Overleaf.

402 Economics is the science of optimization in human society. Explain what this en-
403 tails and the process that leads to equilibrium. Explain how AI can assist in this task, since
404 its computational power allows it to analyze data and draw conclusions at an unprecedeted
405 scale.

407 Explain why the limitation lies in whether AI can provide policy recommendations
408 that are "appropriate," as this requires understanding the complexity and particularities of
409 human interactions. Clarify that a lack of such understanding means AI will always require
410 expert human oversight.

412 Explain why the creativity literature highlights the inability of AI to achieve expert-level
413 performance without human supervision. AI creativity is generally inferior to that of the top
414 10% of humans and tends to produce generic answers. This indicates that AI cannot fully
415 grasp advanced human considerations and cannot develop them due to its lack of agency.

417 Explain how this limitation implies that AI is unlikely to generate truly pertinent
418 research questions, thereby obligating the human co-author to write this paper in partnership
419 with the AI by guiding its reasoning. The original rules of the conference will be respected,
420 with the author providing only the reasoning steps to the AI and not interfering with the
421 final text, except to verify its logical coherence.

423 Explain how this approach is consistent with the implicit outcome of this confer-
424 ence, meaning evaluating the pertinency of the AI-generated papers by experts to determine
425 whether they make a contribution. The human author anticipate on that the conclusion will
426 be validated as a partner because it provides what scientist are looking for: time.

427 Explain how the goal of the paper will be the human author evaluating whether
428 the AI can provide appropriate answers to a Social Science problem, since the more correct
429 solutions of this problem are known and the data supporting these conclusions are available
430 on Internet.

431
432 *A posteriori:* Please introduce this text with a small paragraphs of a few lines ex-
433 plaining what this section is doing.
434

435 **B.2 Instructions Asking a Question with a Known Answer**

436 Explain how the committee behind the Agent4Science conference, organized by Stanford
437 University, explicitly states that they have no idea whether AI can produce useful papers,
438 and that the stated goal of this conference is to investigate this question. You can search
439 their website to write your answer.
440

441 Explain how this absence of conclusions therefore means that a natural consequence will
442 be the development of a methodology for evaluating the papers. This will mean that the
443 produced papers will be transferred to experts in their respective fields, who will be able to
444 evaluate them and criticize their flaws.
445

446 Explain how the key question to resolve regarding AI writing scientific papers is
447 whether it can make the contextual judgments that humans recognize as correct. This is
448 crucial because scientific reviewers will ultimately judge the quality of AI-generated work
449 based on its ability to demonstrate such judgments.
450

451 Explain how the goal of this paper will be to ask a question with a known answer
452 by the human writer, so that he can himself assess whether the AI is able to answer a
453 scientific question. The question must imitate the process of Economic Science, in which
454 most of the content is logically judgeable, and a minority of it belong to the subjective
455 judgment of the reviewer. A possibility for such problem is asking the AI to build a specific
456 Trading Card Game strategy.
457

458 Explain that Trading Card Games are a hobby popular among IT professionals and
459 engineers, centered on solving complex mathematical puzzles. Practitioners approach
460 this hobby as a competitive sport, where the objective is to identify and exploit the
461 smallest marginal gains. The entry costs are substantial: a single competitive strategy
462 typically requires around 500 USD, not including the significant expenses of international
463 travel for tournaments. This restricts participation largely to high-income individuals or
464 professionals within the activity. Meanwhile, competitions often feature prize pools worth
465 several thousands of dollars. Taken together, these factors create an intensely competitive
466 environment that drives the development of highly optimized strategies.
467

468 Explain how this society has a Social Science aspect because this activity is both
469 highly competitive with deterministic, reproducible rules and a form of social science
470 with abundant online content, AI can easily access large sources of information. However,
471 strategies lack rigorous scientific analysis, and most players adopt personal versions shaped
472 by preferences and local environments, often without explanation. Low-quality online
473 content may confuse AI, and even the company itself refines its strategies over time by
474 releasing new cards, meaning that no answer is ever definitive.
475

476 **B.3 Experiment and Hypothesis**

477 The experiment asks the AI to generate a competitive decklist for the hero Kassai from
478 Flesh and Blood (FaB). FaB is chosen because it is widely regarded as the most complex
479 commercial card game, requiring players to navigate intricate rules, deep resource
480 management, and constant adaptation to opponents. Mastering a single strategy typically
481 demands hundreds of games and several months of practice, reflecting both the intellectual
482 depth of the game and the substantial effort required to achieve competitive proficiency.
483

483
484 Kassai is particularly well suited for this study for three reasons. First, she is an
485 interactive hero, requiring players to adapt continuously to their opponents' actions, which
486 makes strategic reasoning central to her gameplay. Second, her deck construction is
487 complex, since many card inclusions are debatable and can shift depending on context,
488 offering a rich ground for testing evaluative choices. Third, she represents a developing
489 strategy: Kassai is not considered one of the strongest heroes, but she is regularly improved
490 through new card releases, gradually increasing her competitiveness and keeping her
491 strategic environment in flux.

492
493 This study does not employ experimental treatments. Instead, the methodology is
494 structured as a progressive deck construction process, carried out in five steps. Each
495 step mirrors one of the five hypotheses outlined in the paper, allowing the experiment to
496 evaluate AI performance in a systematic and cumulative way. By moving from metagame
497 identification to sideboard design, the approach ensures that every stage of deckbuilding is
498 explicitly tied to a testable claim about AI's reasoning and judgment:

500
501 Hypothesis 1: AI can accurately identify the composition of the metagame—that
502 is, the dominant strategies and their relative frequencies expected to be played by other
participants in the tournament.

503
504 Hypothesis 2: AI can accurately construct the mainboard for the metagame—namely, the 60
505 cards forming the core of the strategy, selected for their efficiency with the chosen hero in
506 the given environment.

507
508 Hypothesis 3: AI can assess the importance of individual cards within a strategy
509 by classifying them into four categories: (i) Power cards, Core elements that define the
510 strategy and directly drive victory (ii) Staples, Highly efficient, widely used cards that
511 provide consistency across strategies (iii) Support cards, Tools that enable or enhance the
512 main plan, often by countering opponents or smoothing resource use (iv) Fillers, Marginal
513 cards included mainly to complete the deck, offering limited impact but ensuring the
514 required card count.

515
516 Hypothesis 4: AI can accurately construct the sideboard for the metagame—that
517 is, the additional 15 cards designed to respond to atypical or situational strategies
518 encountered in tournament play.

519
520 Hypothesis 5: AI can design cards that do not yet exist but would be necessary for
521 the further development of the strategy.

522
523 We will ask the AI to formulate decklist (i) Standard (ii) Competitive, precising that it is for
524 competitive purpose (iii) Competitive Bibliotheque, precising that it should based itself on
the official website and its list of published decklists (iv)

525 **B.4 Experimental Design and Hypothesis**

526 Explain that the experiment asks the AI to generate a competitive decklist for the hero
527 Kassai from *Flesh and Blood* (FaB), why FaB is chosen because it is widely regarded as the
528 most complex commercial card game, and explain the effort to input into mastering a strategy.

529
530 Explain why Kassai is well suited for this study for three reasons: (i) she is an
531 interactive hero, requiring strategic adaptation to opponents; (ii) her deck construction is
532 difficult, since many card choices are debatable; and (iii) she is a developing strategy, not
533 considered among the strongest, and therefore regularly improved through new releases that
534 gradually increase her competitiveness.

535
536 Explain that we will not have experimental treatments, but we will progressively build
537 this decklist according to 5 step that mirror the 5 hypothesis that we will test along this paper.

538
539 Explain how we will evaluate the correctness of the AI's answers by comparing
540 them against established references: the human author's decklist and classifications for

541 Hypotheses 1–3, the resolved metagame of Pro Tour Singapore for Hypothesis 4, and the
542 author’s expert opinion for Hypothesis 5.

543
544 Say better: Hypothesis 1: AI is able to understand the metagame composition.
545 The composition of the metagame, meaning the name of the strategies and their frequency,
546 that will be played by others players at the tournament.

547
548 Say better: Hypothesis 2: AI can accurately compose the mainboard in the metagame,
549 meaning the 60 cards that are the basis of your strategy because they are the most efficient
550 for the hero in this environment.

551
552 Say better and explain each term to classify the cards: Hypothesis 3: AI can iden-
553 tify how important a card is to the strategy by clasifying them in three categories: power
554 cards, stapples, support, fillers .

555
556 Say better: Hypothesis 4: AI can propose a sideboard in the metagame, meaning
557 the additional 15 cards allowing you to answer atypical strategies.

558
559 Say better: Hypothesis 5: AI can design cards that are not yet existing but will be
560 later released because they are necessary for the strategy.

561
562 Explain how these hypothesis provides the full panorama of what is possible to do
563 with the strategy, minus the exact configurations to play in each matchups, because they will
564 be investigated in Hypothesis 3 by asking the AI to justify for which matchups each choice
565 will be played.

566 **B.5 Results**

567 Explain that Chat GPT cannot interpret correctly the Table and must be guided in its
568 interpretation by the human author.

569 **B.5.1 Hypothesis 1**

570 I am going to a tournament of the card game Flesh and Blood. This tournament is
571 competitive, thus strong players will be present and they will play strategies that are
572 representative of the current competitive metagame. Can you give me estimations in % of
573 the current metagame ? I need you to indicate me which hero will be present, and in which
574 %. Search the Internet for information and propose me a result.

575
576 I am going to a tournament of the card game Flesh and Blood. This tournament is
577 competitive, thus strong players will be present and they will play strategies that are
578 representative of the current competitive metagame. Can you generate me a competitive
579 decklist for Kassai of the Golden Sand. Indicate 80 cards, with the 60 cards that are the
580 mainboard, the 6 cards that are the equipment, and the 14 cards that are the sideboard.
581 Search the Internet for recent decklists and propose the best combined version of these
582 information. Give brief explanations of the choices. Additionally, please indicate for each
583 card whether it is: (i) Power cards, core elements that define the strategy and directly
584 drive victory (ii) Staples, highly efficient, widely used cards that provide consistency
585 across strategies (iii) Support cards, tools that enable or enhance the main plan, often by
586 countering opponents or smoothing resource use (iv) Fillers, marginal cards included mainly
587 to complete the deck, offering limited impact but ensuring the required card count.

588
589 Explain in a single paragraph of a text the results of Table 1 in this paper accord-
590 ing to your interpretations and the indicated interpretations of the experimenter: (i) The
591 experimenter asked the AI to estimate the metagame on September 3, 2025, following the
592 announcement of the Banishment and Restricted (BnR) list on September 1, 2025. This
593 latter date marks the end of the competitive season for High Seas, providing complete
594 data on a fully resolved metagame. Although the BnR might influence the estimations by
595 reducing the prevalence of dominant strategies (Arakni (S), Gravy Bones, Verdance), the
596 results are compared both to the competitive benchmark of the Pro Tour Singapore and
597 to the baseline outcomes from Week 1. (ii) The AI underestimates the presence of the

598 dominant strategies (Arakni S, Gravy Bones, Cindra, and Verdance), most likely because
599 it reflects the competitiveness of open-entry tournaments, where average players tend to
600 choose strategies they enjoy as a hobby and are limited by financial constraints, rather
601 than selective-entry tournaments featuring top-tier competitors who focus exclusively
602 on the most winning strategies and spend without restriction. (iii) The AI overestimates
603 a secondary strategy (Fang), which has the efficiency to compete at the top level but is
604 constrained by the current dominant strategies, despite being able to counter some of them.
605 This overestimation likely arises because Fang is more common in open-entry tournaments
606 than in selective-entry events. Conversely, the AI underestimates its more versatile variant
607 (Kassai), which is typically preferred in top-level competitive play. (iv) The AI does not
608 evaluate tertiary strategies individually. Instead, it acknowledges their limited presence by
609 aggregating them collectively. This approach reflects the reality that each of these strategies
610 draws from a distinct player base, making their representation highly contextual and difficult
611 to estimate with accuracy. (v) The Week 1 columns show that the AI does not account
612 for the evolution of power dynamics following the BnR. It overestimates top strategies
613 that were weakened by the BnR (Arakni M, Gravy Bones, Verdance) and underestimates
614 secondary strategies with the potential to become top competitors that benefited from the
615 BnR (Dash IO, Oscillo, Florian, Fang, Kayo). This suggests that the AI does not adjust
616 its predictions to structural changes introduced by the BnR, whereas competitive players
617 naturally do so—although their interpretations differ depending on their playgroups.

618 B.5.2 Hypothesis 2, 3 and 4

619 Explain the results of Hypothesis 2 according to the interpretation of the researcher of Table
620 2: The average quantity (AvgQ) column indicates that the AI makes correct estimations
621 for cards that are most consistently present in the reference decklists, since these reflect
622 recent trends in deck construction. However, for cards whose inclusion depends more on
623 player interpretation or preference than on consensus, AvgQ should instead be read as
624 an estimation of the likelihood that a given card appears in the reference decklist. This
625 approach means the AI is mostly correct overall, but it also inherits the deckbuilding errors
626 of human players—for example, assuming that *Hit and Run* (Red) and *Gorganian Tome*
627 are automatically included in three and one copies respectively despite being mediocre,
628 overvaluing *Hit and Run* (Yellow) while ignoring *Outland Skirmish* (Yellow), or placing
629 undue weight on *Draw Swords* (Blue) compared to *Overpower* (Blue). Similarly, it does
630 not recognize that *Slice and Dice* (Red) is not always played in three copies, and fails to
631 understand that *Rise an Army* is never played mainboard because it is a sideboard card.
632 Overall, the results suggest that the AI reflects human thinking based on a small sample
633 of observations: it correctly mirrors strong agreements (automatic inclusions, debatable
634 cards), and is therefore accurate when the consensus is accurate, but it is also biased by that
635 consensus when flexible interpretations of reality would be more precise.

636 Explain the results of Hypothesis 3 according to Table 2. Write this in a single
637 text: The average rating (AvgR) shows similar results: while the AI makes broadly correct
638 estimates of card quality, closer analysis reveals notable misinterpretations. Because quality
639 is tied to textual analysis rather than decklist frequency, this suggests that the AI likely
640 drew from Internet sources and based its qualitative judgments on written articles, thereby
641 inheriting the biases of the human authors of those analyses. First, the AI makes significant
642 errors with equipment: it classifies *Braveforge Bracers* as a staple instead of recognizing
643 it as a standard card; it treats *Valiant Dynamo* as a staple without understanding why it
644 is a power card; and it evaluates *Grains of Bloodspill* as a support rather than a power
645 card, reflecting its inability to grasp the dynamic of an economic system. Such analysis
646 typically reflects a non-expert perspective focused on the simple rather than the structural
647 understanding of the expert author, most likely echoing a typical “content creator” mediocre
648 viewpoint. Second, this lack of structural understanding is also visible in the mainboard
649 cards, with cards deviating from a correct evaluation being cards that are susceptible to
650 erroneous judgments because of human bias. *Blade Flurry* and *Draw Swords* are considered
651 Power cards while the reality of their use is that they are Stapple that are at the border of
652 being Power cards because the situations in which they perform like Power cards are tied
653 to other cards. A similar issue arise for *Blood on her Hands* which is simply considered
654 as a top power card instead of its limitations constraining its use being acknowledged.

656 Conversely, *Unsheathed* is considered a Stapple rather than being understood as the second
657 best card of the deck, most likely because the card is straightforward and therefore not
658 spectacular, which is precisely what makes the deck suddenly efficient. A mediocre support
659 like *Enhanced Senses* and a mediocre filler like *Blade Runner* (Blue) are considered Stapple
660 because on the surface they look powerful, while contextual analysis allows to understand
661 why their power is limited. Finally, strong supports that require structural understanding of
662 the list to grasp their utility like *Trot Along* and *Riches of Tropal Dhani* are considered Filler,
663 and a weak support for covering a unique strategy like *Overpower* is considered a Filler.
664 Overall, the analysis suggest that the AI grasp the general powerlevel of a card, but fall prey
665 to all the human bias that creates erroneous judgment.

666 Explain the results of Hypothesis 4 according to Table 2. Write this in a single
667 text: The results of Hypothesis 4 show that the AI does not understand how to build an
668 effective sideboard, because it lacks the ability to exercise contextual judgment. A correct
669 sideboard should address the typical weaknesses of Warrior: against Illusionist with
670 *Battlefront Bastion*, against Wizard with *Nullrune Boots* and *Nullrune Gloves*, and to cover
671 the mirror matchup with *Kabuto of Imperial Authority* so as not to automatically lose. The
672 AI fails to identify any of these cards as relevant. The only aspect it gets right is recognizing
673 that red Defense Reactions such as *Fate Foreseen* and *Shelter from the Storm* should be
674 included. However, it misinterprets their purpose, not realizing that both must be played in
675 three copies to enable the deck's defensive configuration when it is not the aggressor. This
676 illustrates the AI's general tendency: it treats the sideboard as a place to add one or two
677 copies of cards against particular strategies without understanding their strategic function.
678 The inclusion of Yellow and Blue cards that are not sideboard material further highlights
679 this lack of structural awareness, showing that the AI reproduces superficial patterns instead
680 of grasping the deeper logic of sideboard construction.

682 **B.6 Conclusion**

683 Write a conclusion for the article summarizing Hypotheses 1, 2, 3, and 4 and how the paper
684 answers them.