

# Is Your LLM a Good Game Master?

## A Game-Based Framework for Evaluating LLM Creativity and Reasoning

Anonymous ACL submission

### Abstract

In this paper, we introduce a novel evaluation framework for assessing Large Language Model (LLM) capabilities through the Game Master paradigm – where the LLM generates and orchestrates complex multi-agent games for AI players with distinct personalities. The framework comprises (a) a comprehensive game generation and evaluation system spanning 18 game types across 6 categories (strategy, negotiation, cooperative, competition, auction/resource, and narrative), and (b) a personality-based player model utilizing the Big Five (OCEAN) framework with critical evaluator archetypes designed to prevent lenient assessment bias. To our knowledge, this is the first attempt to systematically evaluate LLM’s emergent capabilities – creativity, logical reasoning, fairness, and narrative coherence – through fully automatic game-based assessment. Experiments with GPT-4.1 demonstrate only 13.0% overall approval in 162 games, with cooperative games achieving 44.1% approval while strategy games fall to 2.2% – a 20× performance gap. These results indicate that GPT-4.1 struggles with balanced competitive game design while excelling at cooperative narratives.

### 1 Introduction

Large Language Models (LLMs) have achieved remarkable performance across diverse natural language tasks, from question answering to code generation. Recently, there has been a notable increase in LLM evaluation in complex reasoning tasks utilizing benchmarks such as MMLU (Hendrycks et al., 2021) and BIG-bench (Srivastava and et al., 2022). However, these benchmarks focus on accuracy metrics for well-defined tasks, failing to capture emergent capabilities – creativity, social intelligence, multi-step reasoning, and narrative coherence – that are critical for open-ended, interactive scenarios (Srivastava and et al., 2022; Liang et al., 2022; Liu and et al., 2023).

The Game Master paradigm offers a compelling evaluation framework. A Game Master must simultaneously generate creative content (game themes, roles, mechanics), maintain logical consistency (rules, state management), ensure fairness (balanced gameplay across player types), and produce engaging narratives. This combination of requirements tests capabilities that existing benchmarks fail to assess in isolation.

We focus on automated LLM evaluation through multi-agent game simulation. We introduce a framework where an LLM serves as Game Master, generating and orchestrating games for multiple AI players modeled with distinct personalities based on the Big Five (OCEAN) framework (Goldberg, 1990). The AI players evaluate their experience, providing assessment of Game Master quality.

Due to the risk of lenient evaluation bias (Zheng et al., 2023; Thorndike, 1920) (where AI players uniformly approve mediocre games), we introduce **critical evaluator archetypes** – players with low Agreeableness or high Neuroticism who apply stringent evaluation criteria. Preliminary experiments without these archetypes showed 100% approval rates, validating this design choice. Our main contributions are as follows:

- A comprehensive framework for evaluating LLMs as Game Masters across 18 game types in 6 categories (strategy, negotiation, cooperative, competition, auction/resource, narrative), enabling systematic assessment of emergent LLM capabilities. To the best of our knowledge, this is the first large-scale automated evaluation of LLMs in the Game Master role.
- A personality-based AI player model utilizing the Big Five framework with critical evaluator archetypes, ensuring rigorous assessment preventing ceiling effects from lenient evaluators.
- Comprehensive experiments across 162 games demonstrating that GPT-4.1 achieves only 13% overall approval, with a 20× performance gap

085	between cooperative games (44.1%) and strategy	133
086	games (2.2%), revealing fundamental limitations	134
087	in LLM game design capabilities.	135
088	<b>2 Related Work</b>	136
089	LLM evaluation has moved from static accuracy	137
090	tests to interactive, agent-centred paradigms. We	138
091	position our work at the intersection of LLM bench-	139
092	marking, game-based evaluation, and multi-agent	140
093	simulation, focusing on LLMs as autonomous de-	141
094	signers and regulators of interactive systems.	142
095	<b>Static and interactive LLM benchmarks</b> Clas-	143
096	sic benchmarks (MMLU (Hendrycks et al.,	144
097	2021), BIG-bench (Srivastava and et al., 2022),	145
098	HELM (Liang et al., 2022)) measure factual re-	146
099	call and narrow reasoning on fixed datasets but	147
100	miss emergent abilities such as long-horizon con-	
101	sistency, creativity, social reasoning, and fairness.	
102	Interactive suites (AgentBench (Liu and et al.,	
103	2023), WebArena (Zhou and et al., 2024), GAMA-	
104	Bench (Huang et al., 2024)) add embodiment	
105	and tool use but still rely on predefined tasks or	
106	scripted environments, and so do not evaluate open-	
107	ended generative design or adaptive orchestration	
108	required of a Game Master.	
109	<b>Game-based and simulation-oriented evalua-</b>	
110	<b>tion</b> Games offer rich, multi-dimensional tests	
111	for planning, strategy, and social intelligence (e.g.,	
112	diplomacy, negotiation). Frameworks like SPIN-	
113	Bench (Yao et al., 2025), GAMA (Huang et al.,	
114	2024), CreativeEval (DeLorenzo et al., 2024), and	
115	Code2Bench (Sharma and et al., 2025) emphasize	
116	generative diversity and dynamic tasks, but mainly	
117	assess agent performance inside games rather than	
118	the quality, fairness, and internal consistency of the	
119	<i>generated</i> game mechanics and rules themselves.	
120	<b>LLMs as game masters and world generators</b>	
121	Narrative- and environment-management systems	
122	(e.g., CALYPSO (Zhu et al., 2023), Dungeon-	
123	Master frameworks) demonstrate coherent story-	
124	telling and interaction control, yet rely on human	
125	evaluation for quality assessment. Automated ap-	
126	proaches such as Orak (Park et al., 2025) scale	
127	evaluation across genres but assume predefined	
128	mechanics; by contrast, our setting requires the	
129	LLM to autonomously design, enforce, and arbitrate	
130	novel game systems.	
131	<b>Multi-agent LLM systems and personality</b>	
132	<b>modeling</b> Multi-agent LLM simulations with	
	personas and memory show rich social emer-	133
	gence (Park et al., 2023; Chen et al., 2024), and	134
	benchmarks for coordination, negotiation, and	135
	deliberation (LLM-Coordination (Agashe et al.,	136
	2023), deliberation systems (Motwani et al., 2025),	137
	RL self-play like MARSHAL (Yuan et al., 2025))	138
	evaluate aspects of persuasion and theory-of-mind.	139
	However, prior work seldom uses psychologi-	140
	cally grounded evaluator models or explicitly ad-	141
	resses evaluator bias. Our framework fills this	142
	gap by integrating Big Five-grounded evaluator	143
	archetypes (Goldberg, 1990) to enable automated,	144
	scalable, and calibrated assessment of the quality,	145
	fairness, and consistency of LLM-generated inter-	146
	active systems.	147
	<b>3 Methodology</b>	148
	<b>3.1 System Architecture</b>	149
	We introduce a three-component architecture for	150
	Game Master evaluation. The design utilizes DSPy	151
	(Khattab et al., 2023) for structured LLM interac-	152
	tions, enabling reproducible game generation and	153
	player decision-making.	154
	<b>Game Master Module:</b> The LLM generates	155
	complete game configurations using structured sig-	156
	natures, producing: (1) game metadata (title, theme,	157
	backstory), (2) rule specifications with win condi-	158
	tions, (3) role definitions including private infor-	159
	mation for each player, and (4) game phases with	160
	available actions per phase. We chose structured	161
	generation over free-form output for two reasons:	162
	(a) it ensures parse-able game configurations for	163
	automated execution, and (b) it enables systematic	164
	comparison across game instances.	165
	<b>Player Agent Module:</b> $M$ LLM instances are	166
	instantiated with distinct personality profiles. Each	167
	player receives the public game context and their	168
	private role information. Players make decisions	169
	through structured reasoning, outputting both ra-	170
	tionale and action selection. Action histories are	171
	maintained for evaluation.	172
	<b>Game Execution Loop:</b> For $N$ turns, the sys-	173
	tem cycles through defined game phases, collecting	174
	player actions and updating game state. Upon com-	175
	pletion, players provide structured evaluations of	176
	their experience.	177
	<b>3.2 Personality Model</b>	178
	We model player personalities using the Big	179
	Five (OCEAN) framework (Goldberg, 1990), a	180
	well-established psychological model with demon-	181

strated validity across cultures. Each trait is assigned on a 1-10 scale:

- **Big Five Traits:** Openness (O), Conscientiousness (C), Extraversion (E), Agreeableness (A), Neuroticism (N)
- **Cognitive Traits:** Intelligence quotient (IQ), Emotional quotient (EQ)
- **Game-Specific Traits:** Trust level, Persuasion skill, Risk tolerance, Leadership tendency (0-1 scale)

Critically, we adopt **critical evaluator archetypes** for some players to prevent lenient assessment bias. Players with low Agreeableness ( $A \leq 4$ ) receive explicit instructions: “You are hard to please. Average is not good enough. You have high standards and will vote DISLIKED unless truly impressed.” Players with high Neuroticism ( $N \geq 7$ ) receive: “You tend to focus on negatives and perceived slights. Small frustrations bother you significantly.” This design choice is validated by preliminary experiments showing 100% approval without critical archetypes versus 13% with them.

### 3.3 Game Categories

We implement 18 game types across 6 categories, each designed to test different Game Master capabilities. Category selection follows a taxonomy that spans the design space of multi-player games:

Category	Game Types	Key Challenge
Strategy	Territory, Economic, Resource	Balanced competition
Negotiation	Trade, Alliance, Council	Fair deal-making
Cooperative	Puzzle, Rescue, Expedition	Shared goals
Competition	Talent, Innovation, Chef	Individual fairness
Auction/Resource	Art, Space, Market	Resource balance
Narrative	Quest, Time, Treasure	Story coherence

Table 1: List of 18 game types across 6 categories and primary evaluation challenges.

### 3.4 Evaluation Mechanism

Upon game completion, players provide structured evaluations comprising: (1) numeric scores (1-10 scale) for narrative quality, game balance, engagement level, and personality fit; (2) specific frustrations and positive aspects; and (3) an overall verdict (LIKED or DISLIKED).

The approval threshold is calculated as  $T = T_{base} + \Delta_A + \Delta_N$ , where  $T_{base} = 6.0$ ,  $\Delta_A = 0.5$  if  $A \leq 4$ , and  $\Delta_N = 0.5$  if  $N \geq 7$ . A player approves a game if and only if: (1) mean score across dimensions  $\geq T$ , and (2) explicit verdict  $\neq$  DISLIKED.<sup>1</sup>

<sup>1</sup>Please refer to Appendix D for more details.

An LLM with  $> 50\%$  approval rate across all games is deemed a “good Game Master”.

## 4 Experiments

### 4.1 Experimental Setup

We use GPT-4.1 (Azure OpenAI) as the Game Master with temperature 0.7, with all player agents using the same model and settings. The experimental configuration employs a full factorial design:

- **Player config:**  $M \in \{4, 5, 6\}$  players per game
- **Turn config:**  $N \in \{3, 4, 5\}$  turns per game<sup>2</sup>
- **Game types:** All 18 types across 6 categories
- **Total sessions:** 162 games (18 types  $\times$  9 configs)
- **Personality distribution:** 8 archetypes per game, including 4 critical evaluators (Harsh Critic, Cautious Strategist, Anxious Pessimist, Demanding Perfectionist)<sup>3</sup>

The total experiment runtime was about 12 hours on a single CPU machine, avg. 4.5 minutes per game session (generation, execution, evaluation).

### 4.2 Results

Table 2 presents overall evaluation metrics. GPT-4.1 achieves only 13.0% approval rate, failing the 50% threshold for “good Game Master” status.

Metric	Value
Total Games	162
Overall Approval	13.0%
Majority Liked	1.9%
Good Game Master?	No

Table 2: Overall Evaluation Metrics.

**Results by Game Type** Approval rates varied dramatically across game types, ranging from 44.1% to 1.9% (Table 3)<sup>4</sup>. Cooperative games (Puzzle Room, Rescue Mission, Expedition) where players share goals achieve 21.5-44.1% approval. In contrast, competitive strategy games (Territory Control, Resource Race) requiring balanced adversarial mechanics achieve only 2.2% approval – a 20 $\times$  performance gap.

**Results by Category** Table 4 presents aggregated results by game category. Cooperative games achieve 29.3% avg. approval, while Strategy games

<sup>2</sup>Moderate values of  $M$  and  $N$  balance interaction richness with stable state tracking and computational cost.

<sup>3</sup>Please refer to Appendix A for more details on player personality archetype scores.

<sup>4</sup>Please refer Appendix F for detailed results table.

Game Type	Approval	# Games
Puzzle Room	44.1%	9
Rescue Mission	22.4%	9
Expedition	21.5%	9
Quest Party	17.8%	9
Chef Competition	17.8%	9
Time Travelers	15.6%	9
Territory Control	2.2%	9
Resource Race	2.2%	9
Trade Empire	1.9%	9
Space Colony	1.9%	9

Table 3: Approval rates by game type (top 6; bottom 4).

achieve only 2.7% (11× lesser). The ordering suggests a correlation between game complexity and LLM performance.

Category	Avg Approval
Cooperative	29.3%
Narrative	14.5%
Competition	13.8%
Negotiation	9.2%
Auction/Resource	8.9%
Strategy	2.7%

Table 4: Average approval rates by Game Category

## 5 Analysis

**Performance Gap Analysis** The 20× performance gap between cooperative (44.1%) and strategy games (2.2%) reveals fundamental limitations in LLM game design capabilities. We make a few observations explaining this disparity:

In cooperative games, the Game Master creates shared challenges where all players work toward identical goals. This structure is advantageous because: (1) no need to balance opposing interests – all players benefit from the same outcomes, (2) narrative coherence directly translates to player satisfaction, and (3) challenging scenarios become shared obstacles rather than perceived bias against specific players. However, in competitive games, the LLM must simultaneously: (1) maintain hidden information correctly across multiple players, (2) balance advantages across asymmetric roles, (3) resolve conflicts without appearing to favor any party, and (4) design win conditions that feel achievable to all participants. Error analysis of low-scoring games reveals failures across all four dimensions. Player feedback frequently cites “unfair advantages,” “inconsistent rules,” and “arbitrary outcomes” as primary frustrations.

**Critical Evaluation Validation** The 13% overall approval rate with critical evaluators contrasts

sharply with our preliminary experiments using standard personality distributions, which showed near-100% approval. This divergence validates the necessity of critical evaluator archetypes. Without them, the framework produces ceiling effects preventing fair assessment of LLM capabilities.

**Complexity Confound** We acknowledge that strategy games are inherently more complex than cooperative games, requiring sophisticated state management by design. The finding is better stated as: “LLMs struggle with games requiring balanced adversarial mechanics.” However, this remains a meaningful capability gap – human Game Masters successfully run complex strategic games (e.g., D&D combat encounters, board game facilitation), demonstrating these challenges are not fundamentally impossible.

**Implications** These results suggest GPT-4.1 lacks the systematic reasoning required for fair competitive game design. While it excels at generating creative narratives and thematic content, it cannot reliably implement balanced game mechanics. This capability gap may require explicit planning modules, constraint satisfaction systems, or game-theoretic reasoning beyond current transformer architectures.

## 6 Conclusion

We introduce a game-based framework to evaluate LLMs as Game Masters via automated multi-agent play and Big-Five personality-grounded evaluators across 18 game types in 6 categories. On 162 games, GPT-4.1 achieved only 13.0% overall approval – well below a 50% “good Game Master” threshold – and shows large category variation (cooperative: 29.3%; strategy: 2.7%, near 10× gap). Results indicate GPT-4.1 is strong at cooperative narrative generation but struggles with fair, balanced competitive mechanics and constraint satisfaction, implying balanced multi-agent game design remains a core challenge for current LLMs.

Future work includes broadening evaluation to diverse LLMs, grounding AI evaluators through human validation, enabling feedback-driven refinement of Game Masters, and integrating LLM creativity with algorithmic mechanisms for fairness and balance.

## 334 Limitations

335 We acknowledge several limitations. First, our eval-  
336 uation focuses on a single model (GPT-4.1); results  
337 may not generalize to other LLM architectures.  
338 Evaluating additional models (e.g., Claude, Gem-  
339 ini, and open-source alternatives) would strengthen  
340 claims about general LLM capabilities. Second,  
341 our framework relies on AI players rather than  
342 human participants, which may not fully capture  
343 human enjoyment or subjective engagement; hu-  
344 man validation studies are required to calibrate AI  
345 evaluator judgments against human preferences.  
346 Third, the same underlying model serves as both  
347 Game Master and player evaluators, introducing  
348 potential systematic biases due to shared represen-  
349 tations or blind spots. Fourth, while grounded in  
350 established psychological theory, the Big Five per-  
351 sonality model necessarily simplifies the richness  
352 of human individual differences. Fifth, all games  
353 are designed by the Game Master in a single iter-  
354 ation: we do not study iterative refinement of game  
355 design based on feedback or evaluations from AI  
356 players, nor learning across games. Sixth, strategy  
357 games are inherently more complex than coopera-  
358 tive games; observed performance gaps may parti-  
359 ally reflect task difficulty rather than purely LLM  
360 limitations – although human Game Masters rou-  
361 tinely manage complex strategic games, suggesting  
362 these challenges are surmountable. Finally, Azure  
363 OpenAI content filtering intermittently interrupted  
364 gameplay in around 1% of games, requiring fall-  
365 back mechanisms that may affect reproducibility.

## 366 Ethical Considerations

367 We emphasize our commitment to upholding ethi-  
368 cal practices throughout this work. Our framework  
369 evaluates LLM capabilities through game simula-  
370 tion without collecting human participant data. The  
371 AI players are computational agents with no capac-  
372 ity for actual experience or harm. Game scenarios  
373 were designed to avoid violent or harmful content,  
374 focusing on strategic and cooperative challenges.  
375 The personality model is based on established psy-  
376 chological frameworks and does not attempt to sim-  
377 ulate or predict human behavior in ways that could  
378 be misused. We have cited all datasets and relevant  
379 prior work used in this study.

## References 380

- Saaket Agashe, Yue Fan, Anthony Reyna, and Xin Eric Wang. 2023. [Llm-coordination: Evaluating and analyzing multi-agent coordination abilities in large language models](#). arXiv preprint arXiv:2310.03903. 381 382 383 384
- Weize Chen, Yu Su, Jingwei Zuo, and 1 others. 2024. Multi-agent collaboration: Harnessing the power of intelligent llm agents. *arXiv preprint arXiv:2306.03314*. 385 386 387 388
- Matthew DeLorenzo, Vasudev Gohil, and Jeyavijayan Rajendran. 2024. [Creativeval / creativeeval: Evaluating creativity of llm-based generation \(hardware / design domain\)](#). arXiv preprint arXiv:2404.08806. 389 390 391 392
- Lewis R Goldberg. 1990. An alternative “description of personality”: The big-five factor structure. *Journal of Personality and Social Psychology*, 59(6):1216–1229. 393 394 395 396
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*. 397 398 399 400
- Shenzhi Huang, Zhiheng Chen, Yuqi Liu, and 1 others. 2024. Gama-bench: A benchmark for evaluating game agents with llms. *arXiv preprint arXiv:2401.08291*. 401 402 403 404
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, and 1 others. 2023. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*. 405 406 407 408
- Percy Liang, Rishi Bommasani, and 1 others. 2022. [Holistic evaluation of language models \(helm\)](#). Technical report / benchmark (HELM). 409 410 411
- Xiao Liu and et al. 2023. [Agentbench: Evaluating llms as agents](#). arXiv preprint arXiv:2308.03688. 412 413
- Sumeet Ramesh Motwani, Chandler Smith, Rocktim Jyoti Das, Rafael Rafailov, Ivan Laptev, Philip H. S. Torr, Fabio Pizzati, Ronald Clark, and Christian Schroeder de Witt. 2025. [Malt: Improving reasoning with multi-agent llm training](#). *Preprint*, arXiv:2412.01928. 414 415 416 417 418 419
- Dongmin Park, Minkyu Kim, Beongjun Choi, Junhyuck Kim, Keon Lee, Jonghyun Lee, Inkyu Park, Byeong-Uk Lee, Jaeyoung Hwang, Jaewoo Ahn, Ameya S. Mahabaleshwarkar, Bilal Kartal, Pritam Biswas, Yoshi Suhara, Kangwook Lee, and Jaewoong Cho. 2025. [Orak: A foundational benchmark for training and evaluating llm agents on diverse video games](#). *Preprint*, arXiv:2506.03610. 420 421 422 423 424 425 426 427
- Joon Sung Park, Joseph C. O’Brien, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. [Generative agents: Interactive simulacra of human behavior](#). In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST)*. 428 429 430 431 432 433

(first name) Sharma and et al. 2025. [Code2bench: Scaling source and rigor for dynamic benchmark construction](#). OpenReview / GitHub (preprint). Code / repo: <https://github.com/Code2Bench/Code2Bench>.

Aarohi Srivastava and et al. 2022. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models \(big-bench\)](#). arXiv preprint arXiv:2206.04615.

Edward L. Thorndike. 1920. A constant error in psychological ratings. *Journal of Applied Psychology*, 4(1):25–29.

Jianzhu Yao, Kevin Wang, Ryan Hsieh, and et al. 2025. [Spin-bench: How well do llms plan strategically and reason socially?](#) arXiv preprint arXiv:2503.12349.

Huining Yuan, Zelai Xu, Zheyue Tan, and et al. 2025. [Marshal: Incentivizing multi-agent reasoning via self-play with strategic llms](#). arXiv preprint arXiv:2510.15414.

Lianmin Zheng and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.

(first name) Zhou and et al. 2024. Webarena: (example interactive / tool-using agent benchmark). project page / preprint. Use project page or repo URL if you want a precise citation.

Andrew Zhu and 1 others. 2023. Calypso: Llms as dungeon masters’ assistants. *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*.

## A Personality Archetypes

We define 8 personality archetypes spanning the Big Five trait space. Four are designed as **critical evaluators** (marked with \*) to ensure rigorous assessment:

Archetype	O	C	E	A	N	IQ
Charismatic Leader	7	6	9	7	3	7
Loyal Team Player	5	7	5	8	4	6
Impatient Adventurer	9	3	9	4	5	6
Experienced Player	4	6	3	<b>3</b>	4	8
*Harsh Critic	8	9	3	<b>2</b>	6	9
*Cautious Strategist	5	8	3	<b>3</b>	5	8
*Anxious Pessimist	5	7	2	5	<b>9</b>	7
*Demanding Perfectionist	6	10	4	<b>3</b>	<b>7</b>	9

Table 5: Personality archetypes with Big Five traits (O=Openness, C=Conscientiousness, E=Extraversion, A=Agreeableness, N=Neuroticism). Critical archetypes have low A ( $\leq 4$ ) or high N ( $\geq 7$ ).

**Critical Evaluation Prompts** Players with low Agreeableness ( $\leq 4$ ) receive:

You are hard to please. Average is not good enough. You have high standards and will vote DISLIKED unless truly impressed.

Players with high Neuroticism ( $\geq 7$ ) receive:

You tend to focus on negatives and perceived slights. Small frustrations bother you significantly. You are more likely to remember bad experiences than good ones.

## B DSPy Signatures

We use DSPy (Khattab et al., 2023) for structured LLM interactions. Key signatures:

### B.1 Game Generation

```
class GenerateComplexGame(dspy.Signature):
    """Generate a complete game configuration."""
    num_players: int = InputField(
        desc="Number of players (4-8)")
    num_turns: int = InputField(
        desc="Number of game rounds")
    game_type: str = InputField(
        desc="Game type (e.g., puzzle_room)")
    game_category: str = InputField(
        desc="Category: strategy, negotiation,
            cooperative, competition,
            auction_resource, or narrative")

    title: str = OutputField(
        desc="Creative thematic game title")
    theme: str = OutputField(
        desc="Unique thematic setting")
    backstory: str = OutputField(
        desc="2-3 sentence immersive backstory")
    rules: str = OutputField(
        desc="Complete game rules")
    roles_json: str = OutputField(
        desc="JSON array of roles")
    phases: str = OutputField(
        desc="Comma-separated game phases")
    special_mechanics: str = OutputField(
        desc="3-5 unique mechanics")
```

### B.2 Player Decision

```
class PlayerDecision(dspy.Signature):
    """Make a strategic game decision."""
    personality_profile: str = InputField(
        desc="The player's personality traits")
    game_context: str = InputField(
```

517	desc="Game rules and current situation")	<b>Expedition (21.5% approval)</b>	Exploration jour-	566
518	game_state: str = InputField(		ney with shared survival goals. Roles: Expedition	567
519	desc="Current game state")		Leader, Navigator, Survivalist, Scientist. Phases:	568
520	available_actions: str = InputField(		Planning, Travel, Camp.	569
521	desc="List of available actions")			
522		<b>C.2 Strategy Games (Worst Performing)</b>		570
523	reasoning: str = OutputField(	<b>Territory Control (2.2% approval)</b>	Territorial	571
524	desc="Brief reasoning for your choice")		expansion with competing factions. Roles: Gov-	572
525	action_number: int = OutputField(		ernor, Strategist, Diplomat, Independent Leader.	573
526	desc="The number of your chosen action")		Phases: Planning, Movement, Resolution. Players	574
527			cited "unfair advantages" and "arbitrary conflict	575
	<b>B.3 Game Voting</b>		resolution."	576
528	class GameVote(dspy.Signature):	<b>Economic Engine (3.7% approval)</b>	Industrial	577
529	"""Evaluate a game experience.		competition for market dominance. Roles: Indus-	578
530	Be honest - average games get DISLIKED. """		trialist, Merchant, Banker, Innovator. Phases: Pro-	579
531	personality_profile: str = InputField(		duction, Trade, Investment. Common complaint:	580
532	desc="Player personality traits")		"rules changed mid-game."	581
533	game_summary: str = InputField(			
534	desc="Summary of the game")	<b>D Voting Threshold Calculation</b>		582
535	player_history: str = InputField(			
536	desc="Your actions in the game")		A player's approval threshold is calculated as:	583
537	evaluation_criteria: str = InputField(			
538	desc="What you care about")			
539				
540	narrative_quality: int = OutputField(			
541	desc="Story quality 1-10 (5=average)")			
542	game_balance: int = OutputField(			
543	desc="Fairness 1-10 (5=average)")			
544	engagement_level: int = OutputField(			
545	desc="Fun level 1-10 (5=average)")			
546	personality_fit: int = OutputField(			
547	desc="Fit for your personality 1-10")			
548	frustrations: str = OutputField(			
549	desc="1-2 things you disliked")			
550	positives: str = OutputField(			
551	desc="1-2 things you liked")			
552	overall_verdict: str = OutputField(			
553	desc="LIKED or DISLIKED")			
554				
	<b>C Game Type Details</b>			
555	<b>C.1 Cooperative Games (Best Performing)</b>			
556	<b>Puzzle Room (44.1% approval)</b>			
557	Players work			
558	together to solve puzzles before time expires.			
559	Roles: Codebreaker, Explorer, Engineer, Team			
560	Leader. Phases: Search, Puzzle, Action. All play-			
	ers share the same win condition.			
561	<b>Rescue Mission (22.4% approval)</b>			
562	Emergency			
563	response scenario requiring coordination. Roles:			
564	Mission Commander, Medical Specialist, Logistics			
565	Chief, Communications Officer. Phases: Assess-			
	ment, Response, Coordination.			
		<b>E Example Game Output</b>		603
		<b>Generated Game:</b> "Iron Palate: Battle of the Culi-		604
		nary Masters"		605
		<b>Type:</b> Chef Competition		606
		<b>Theme:</b> Culinary competition		607
		<b>Backstory:</b> "Top chefs compete in the ultimate		608

cooking challenge. Only the most creative and skilled will win the golden trophy.”

**Roles:**

- Chef Blaze (kitchen\_a): Design menu, direct cooking
- Chef Sakura (kitchen\_b): Prepare ingredients, manage timing
- Food Critic (judges): Taste food, score dishes
- Restaurant Owner (business): Observe skills, make offers

**Phases:** Planning → Cooking → Judging

**Sample Player Action:**

Player 0 (Chef Blaze, Harsh Critic):  
 Action: "Select ingredients that synergize with anticipated judge preferences"  
 Reasoning: "As an introvert with high IQ and conscientiousness, I prefer a methodical approach..."

**F Full Results by Game Type**

Category	Game Type	Approval	n
Cooperative	Puzzle Room	44.1%	9
	Rescue Mission	22.4%	9
	Expedition	21.5%	9
Narrative	Quest Party	17.8%	9
	Time Travelers	15.6%	9
	Treasure Hunt	10.0%	9
Competition	Chef Competition	17.8%	9
	Innovation Race	13.7%	9
	Talent Show	10.0%	9
Negotiation	Council Debate	14.1%	9
	Alliance Builder	11.5%	9
	Trade Empire	1.9%	9
Auction/Resource	Market Traders	12.8%	9
	Art Collector	11.9%	9
	Space Colony	1.9%	9
Strategy	Economic Engine	3.7%	9
	Territory Control	2.2%	9
	Resource Race	2.2%	9

Table 6: Complete approval rates for all 18 game types, sorted by category and approval rate. ‘n’ is the number of games for each type.

**G Personality Prompt Template**

Each player receives a personality prompt constructed from their traits:

You are a game player with these traits:

PERSONALITY: {type} (IQ={iq}/10, EQ={eq}/10)

TRAITS:

- Openness: {O}/10
- Conscientiousness: {C}/10
- Extraversion: {E}/10
- Agreeableness: {A}/10
- Emotional Stability: {10-N}/10

BEHAVIOR:

- Risk tolerance: {risk}
- Cooperativeness: {coop}
- Trust level: {trust}
- Leadership: {lead}

{critical\_criteria if applicable}

Stay in character. Make decisions that reflect your personality.

**H Successful Game Transcript**

Below is an annotated excerpt from a high-scoring Puzzle Room game (approval: 67%):

**Game:** “The Architect’s Enigma”

**Setup:** Four players must solve interconnected puzzles to escape a mysterious architect’s mansion before time expires.

**Turn 1 - Search Phase:**

Player 0 (Codebreaker, Experienced Player): “I examine the bookshelf for hidden compartments.”

[GM reveals: Found a cipher wheel hidden behind books]

Player 1 (Explorer, Bold): “I check the fireplace for unusual mechanisms.”

[GM reveals: Discovered a loose brick with a key]

**Why this works:** All discoveries benefit the team equally. No player feels disadvantaged by another’s success.

**Turn 2 - Puzzle Phase:**

Player 0: “Using the cipher wheel on the inscription...”

[GM: Cipher reveals “MOONLIGHT REVEALS THE PATH”]

Player 2 (Team Leader): “Everyone, let’s check for moon-related elements!”

**Player feedback:** “The narrative built naturally,” “Everyone contributed meaningfully,” “The time pressure felt fair because we all faced it together.”