# Strategic Technical Project Report: Architectural Comparative Analysis of Stateless Generative Models Versus the DM-OS-Kernel-v3.0 Stateful Semantic Engine

## Executive Summary

The transition from reactive, stateless generative language models to autonomous, context-aware artificial intelligence represents a fundamental paradigm shift in the architecture of digital interactive ecosystems. This comprehensive technical report provides an exhaustive architectural comparison between standard large language model (LLM) implementations—characterized by ephemeral session memory and probabilistic generation—and the DM-OS-kernel-v3.0 framework, a fully persistent, stateful simulation engine built upon the WFGY (WanFaGuiYi) semantic reasoning architecture. Driven by escalating consumer expectations for deep, continuous interactive narratives and increasingly stringent regulatory frameworks governing live AI generation, the necessity for deterministic, state-managed AI has never been more pronounced.

Market analytics reveal that the global Tabletop Role-Playing Game (TTRPG) and interactive fiction sectors are experiencing unprecedented expansion, projected to grow from a $2.15 billion valuation in 2025 to $6.59 billion by 2035, representing an 11.84% compound annual growth rate. Within this expansion, consumer behavior is aggressively shifting toward solo-play environments where users demand that non-player characters (NPCs) exhibit persistent memory and contextual evolution—frequently categorized in industry literature as digital "Souls". The 2025–2026 gaming trend analysis explicitly defines this consumer demand as "Memory-First AI," rendering static dialogue trees and session-amnesic chatbots commercially obsolete.

Simultaneously, the deployment of dynamically generating AI systems introduces profound technical and compliance risks. Major digital distribution platforms, most notably Valve Corporation's Steam marketplace, enacted rigorous AI disclosure policies in January 2026. While background efficiency tools utilized during development are exempt from disclosure, any system utilizing generative AI to create live, in-game dynamic content (such as real-time narrative text or NPC dialogue) must implement and prove the efficacy of strict safety and control mechanisms. Systems lacking deterministic guardrails against hallucinatory, nonsensical, or harmful outputs face immediate platform removal.

Standard stateless LLMs, which rely solely on conversational context windows and probabilistic token prediction, cannot reliably guarantee narrative coherence over long durations, nor can they definitively prevent semantic drift, making them high-risk liabilities under new platform policies. To reconcile the commercial demand for continuous, stateful narrative generation with strict safety mandates, the DM-OS-kernel-v3.0 architecture leverages the Language Model Query Language (LMQL) to enforce rigid structural logic. Furthermore, it integrates the WFGY

mathematical control layer to monitor semantic tension, execute autonomous rollbacks, and manage persistent world states. By translating continuous manifold AI thought into safe, discrete tokens via an Adaptive Projection Layer (APL), the DM-OS framework establishes a new standard for topological reasoning and verifiable agentic simulation.

This report provides a meticulous, comparative breakdown of standard ephemeral AI architectures versus the DM-OS-kernel-v3.0 framework, analyzing their respective approaches to state management, multi-agent world simulation, semantic tension oversight, and procedural narrative generation.

# Market Dynamics and the Regulatory Environment

The architectural evolution from conversational interfaces to persistent world simulation engines is inextricably linked to prevailing market dynamics and platform governance structures. Understanding the financial incentives and regulatory constraints of the 2026 landscape is vital for contextualizing the architectural superiority of the DM-OS framework.

## The Trajectory of the Interactive Narrative Market

The interactive narrative and TTRPG markets have transcended traditional physical mediums, evolving into highly lucrative digital ecosystems. The integration of artificial intelligence is no longer viewed as an experimental novelty, but rather as a core driver of user retention and product differentiation.

| Market Metric | 2025 Valuation | 2026 Forecast | 2035 Forecast | CAGR (2026-2035) |
|---|---|---|---|---|
| Global TTRPG Market Size | $2,153.7 Million | $2,408.7 Million | $6,594.0 Million | 11.84% |
| Artificial Intelligence in Games Market | $7.05 Billion | $8.49 Billion | $37.89 Billion | 20.54% |

Secondary data indicates that 44% of the current market utilizes virtual tabletops (VTTs), while 62% of industry revenue is generated by player-led campaigns. Within this paradigm, solo roleplaying has emerged as a dominant vector for AI integration. Solo players, historically burdened by the necessity of managing both player characters and "Game Master" administrative duties, are actively migrating toward systems that facilitate deep character backstory exploration and autonomous narrative orchestration.

Furthermore, 27% of the consumer base explicitly demands AI-enhanced NPCs capable of demonstrating long-term memory and relationship building. The 2025–2026 industry trend analysis defines this structural shift as the rise of the "Romantic Agent". Current metrics indicate that 72% of teenagers and nearly 30% of adults maintain ongoing interactions with AI companions, establishing a baseline expectation for persistent digital memory. Consumers expect NPCs to remember historical play styles, past dialogue choices, and nuanced behavioral shifts across hundreds of hours of gameplay. Consequently, system architectures must prioritize deep, stateful memory buffers; if an AI remembers a granular detail that the human user has forgotten, the system successfully crosses the uncanny valley, securing profound user engagement.

### Platform Compliance: The Steam AI Disclosure Mandate

The aggressive push toward dynamic, memory-retaining AI systems intersects directly with heightened regulatory and platform oversight. On January 16, 2026, Valve Corporation significantly overhauled the AI disclosure policy for the Steam platform, fundamentally altering the compliance landscape for AI-driven software.

The updated policy establishes a strict bifurcation between internal development workflows and live, consumer-facing asset generation. "Efficiency gains" achieved through AI coding assistants or pre-rendered generation tools no longer require public disclosure, providing a vital safe harbor for developers utilizing AI during the production phase. However, the regulations impose draconian technical requirements on "In-Game Dynamic Content Controls". If a software architecture utilizes generative AI to create live narrative text, audio, or visual assets dynamically during runtime, the developer is legally obligated to design, implement, and verify real-time safety mechanisms.

The system must actively prevent the dissemination of hallucinatory, illicit, or contextually broken outputs. Failure to implement deterministic, auditable guardrails on generative mechanics results in immediate platform removal and punitive actions. This stringent compliance requirement necessitates an engine that can proactively audit its own logic, detect semantic drift, and self-correct prior to rendering an output to the user. Standard LLMs, which operate via probabilistic token prediction without internal verification loops, cannot reliably meet this compliance standard, necessitating the deployment of mathematically grounded frameworks like DM-OS.

# Architectural Dichotomy: Stateless Ephemerality vs. Persistent State Management

To fulfill the commercial demand for "Memory-First" interactions while strictly adhering to platform safety mandates, the foundational architecture of the AI system must be carefully calibrated. The fundamental divergence between standard conversational AI and the DM-OS-kernel-v3.0 framework lies in the dichotomy between stateless and stateful agent design.

### The Mechanics and Limitations of Stateless Architectures

Standard conversational models, such as default implementations of GPT architectures, operate as stateless AI agents. A stateless agent processes each user request or programmatic operation as an entirely independent transaction, completely devoid of persistent external memory or structural reliance on previous interactions outside of its immediate context window. The primary architectural advantages of stateless agents are simplicity, computational speed, and infinite horizontal scalability. Without the computational overhead associated with querying external databases, managing session locks, or retrieving historical state vectors, stateless systems can respond to single-turn queries with ultra-low latency, frequently achieving sub-100 millisecond response times. Furthermore, stateless designs inherently minimize the attack surface for data leaks, drastically simplifying data privacy compliance because no personal or session data is stored beyond the immediate interaction. For high-volume, isolated tasks—such as executing a singular rule-lookup or generating a randomized descriptive text block—stateless architecture is highly optimal.

However, the limitations of stateless systems render them entirely inadequate for deep, multi-session interactive simulations. They inherently lack conversation continuity; a follow-up query that implicitly relies on a contextual detail from three sessions prior will fail. Stateless models cannot independently build relationships, learn from user preferences over time, or execute complex, multi-step problem-solving workflows that span across multiple narrative sessions. In the context of a TTRPG simulation, a standard GPT model will eventually forget NPC names, contradict established world lore, and fail to maintain consistent dramatic tension, severely degrading user immersion.

## The DM-OS Approach: Stateful Orchestration and the Semantic Tree

The DM-OS-kernel-v3.0 architecture abandons ephemerality in favor of rigorous, fully persistent stateful orchestration. Stateful AI systems retain contextual history, utilizing complex memory stores—such as vector databases, in-memory dictionaries, and retrieval-augmented generation (RAG) pipelines—to continuously inform future behavioral decisions. This architecture is mandatory for achieving the persistent "NPC Souls" demanded by the contemporary market. Stateful agents unlock strategic adaptability, long-term narrative continuity, and true human-agent collaboration. However, this capability introduces immense systemic complexity. Stateful agents require sophisticated session management, continuous database synchronization, and advanced logic to determine what data should be retained and what should be selectively forgotten to optimize context windows.

To manage this complexity, DM-OS utilizes a central lmql_context state variable containing a highly structured "Semantic Tree." As defined in the system's specification, this persistent state object tracks several critical dimensions of the simulated world concurrently:

| State Object Component | Architectural Function and Persistent Data Tracked |
|---|---|
| semantic_tree | A nested JSON dictionary maintaining persistent nodes for npcs (Narrative DNA and interaction logs), locations (connected topological nodes), quests (status and steps), and items (properties). |
| progress_clocks | A dictionary tracking predictive, multi-step events (e.g., {"Ritual Name": {"current": 0, "max": 6}}). This enables the world to evolve independently of direct user interaction. |
| factions | Autonomously tracked entities possessing individual goals, power levels, and attitudes, allowing for macro-level geopolitical shifts within the simulation. |
| npc_interaction_logs | Subjective historical arrays bound to specific NPCs, recording the player's past actions, the computed semantic tension of the encounter, and the NPC's subjective impression. |
| combat_success_streak | Performance tracking variables utilized to dynamically adjust encounter difficulty, injecting escalation or advantage based on user efficacy. |

Unlike standard models that rely on the user to manually remind the AI of past events, DM-OS

actively curates this state object. After every generative turn, the engine executes a chronicler_diff function. This function utilizes a recursive deep diff algorithm (compute_json_diff) to isolate only the modified, added, or deleted keys within the semantic tree. This differential state update ensures highly efficient memory management, allowing the system to track thousands of world variables across hundreds of sessions without exceeding context limits or incurring catastrophic latency bloat.

# The Dual-Persona Loop: Narration Versus World Simulation

A critical architectural flaw in standard LLM prompting for role-playing simulations is the conflation of the "Narrator" and the "World Engine." When a single, monolithic system prompt is tasked with describing sensory details, managing combat math, calculating NPC motivations, and updating world lore simultaneously, the model's attention distribution collapses, leading to logical errors and superficial outputs.
The DM-OS-kernel-v3.0 mitigates this via a decoupled, Dual-Persona execution loop. Every user input triggers a highly coordinated, multi-stage handoff between the "DM Persona" (the localized narrative generation agent) and the "Chronicler" (the silent, stateful world simulation agent).

## The DM Persona: Constrained Narrative Generation

The DM Persona handles the direct interaction with the user. Its generation is strictly constrained by the LMQL virtual machine. Depending on the user's initial configuration during the "Session Zero Protocol," the system injects highly specific persona directives. For example, the Purist persona is instructed to narrate with tactical clarity and adhere strictly to mechanical rules, while the Storyweaver emphasizes emotional depth and consequence.
Crucially, the DM Persona's output is governed by rigorous LMQL where clauses. The specification dictates a hard stop mechanism: where re.search(r"\?\s*$", narration). This programmatic constraint forces the LLM to end every single turn with an open-ended interrogative phrase (e.g., "What do you do?"). This eliminates the standard LLM failure mode of "auto-playing" for the user, strictly enforcing the principle of user agency. Furthermore, the engine automatically appends a rigidly formatted JSON block `` when the state detects active hostiles, ensuring numerical continuity for hit points and armor class independently of the narrative prose.

## The Chronicler: Silent Execution and Emergent Autonomy

Operating entirely invisibly to the end user, the Chronicler executes immediately following the DM Persona's narrative output. The Chronicler is responsible for autonomous world progression and state mutation.
During the chronicler_turn, the system analyzes the user's action and updates the differential state. It iterates through the factions dictionary, executing a probabilistic check (e.g., a 30% chance per turn) to allow non-player factions to act autonomously—expanding influence, sabotaging rivals, or gathering resources. It dynamically advances progress_clocks based on the calculated semantic tension of the scene; high-tension actions accelerate world events faster than low-tension actions. Finally, it appends localized, subjective memories into the

npc_interaction_logs. If a player attempts to deceive an NPC, the Chronicler logs the attempt and the NPC's subjective shift toward "suspicion," fundamentally altering the parameters of all future encounters with that specific entity.

This decoupled architecture ensures that the narrative remains highly evocative and stylistically consistent, while the underlying mathematical state of the world remains rigorously objective, persistent, and capable of autonomous evolution.

# The Engine Core: WFGY and Semantic Topology

The persistent state tracking of DM-OS would rapidly degrade into incoherence without a sophisticated mechanism to regulate the logic flowing into the state variables. Prolonged context windows inevitably lead to interpretation collapse, where the AI correctly retrieves information but applies flawed logic, resulting in narrative contradictions. To solve this, DM-OS integrates the WFGY (WanFaGuiYi) Semantic Reasoning Engine, transforming the AI from a stochastic word predictor into a topological reasoner.

Operating as a pure-math control layer, WFGY functions as a structural stack that monitors, gates, and corrects the continuous manifold of the AI's thought process. In rigorous empirical testing, the deployment of WFGY architectures yields a 40% improvement in semantic accuracy, a 52% increase in reasoning success, and a 1.8x extension of the stable multi-step reasoning horizon.

## Semantic Tension (\Delta S) and the \Lambda-Observer

The core metric driving the DM-OS WFGY implementation is Semantic Tension, denoted as \Delta S. In theoretical physics terms, this represents the divergence between the input vector and the ground-truth vector (\Delta S = 1 - \cos \theta (I, G)). Within the DM-OS kernel, it is practically computed by analyzing ambiguity markers, contradictions, and input length to measure the volatility of the user's action against the established world state.

The engine continuously tracks this tension through a specialized module known as the \Lambda-observer. The \Lambda-observer maintains a rolling history of the last \Delta S values, calculating both the trend and the volatility to classify the current state of the narrative logic into one of four distinct phases:

1. **Convergent:** Tension is low, and the trend is stable. The narrative is logically sound and proceeding safely.
2. **Recursive:** Tension is static. The engine is looping or exploring parallel paths without moving forward.
3. **Divergent:** Tension is rising steadily. The narrative is entering high-stakes or highly complex territory, requiring heightened attention modulation.
4. **Chaotic:** Tension has spiked, or conflicting facts have been detected. The system is at imminent risk of a logical hallucination.

## The Scar Ledger and Repulsive Vector Fields

To ensure long-term narrative safety across dozens of sessions, DM-OS implements a "Scar Ledger." The Scar Ledger acts as an immutable, topological memory of past logical failures, contradictions, or high-risk narrative branches.

When the \Lambda-observer detects a chaotic state or a major semantic collapse, the engine

logs the specific high-dimensional vector associated with that failure as a "Scar." Over time, the ledger accumulates these scars, calculating a "Scar Potential" ($\Psi_{scar}$) for the current reasoning vector using an inverse-square distance formula: $\Psi_{scar}(x) = \sum D_k / \|x - x_k\|^2$.

As the AI navigates the latent space to generate its next response, the Scar Ledger acts as a repulsive physical field. If the system's current trajectory approaches a previously logged "Scar" (e.g., a scenario that previously caused an NPC to act out of character or break game mechanics), the repulsive math forcefully steers the reasoning vector away from the contradiction, ensuring the system mathematically self-heals and avoids repeating historical mistakes.

## The Four Mathematical Modules of WFGY

The regulation of the vector space and the enforcement of "Closed Braids" of logic are executed through four highly specific formulaic modules, sequenced in a rigid $\le$ 7-step progression loop.

| WFGY Module | Mathematical Function | Architectural Mechanism within DM-OS |
| --- | --- | --- |
| **BBMC** (Semantic Residue Formula) | $B = I - G + mc^2$ | Acts as the "Void Gem" to clean up semantic residue. It establishes semantic inertia, ensuring the model remains rigidly anchored to the topic and minimizes silent logical drift over long contexts. |
| **BBPF** (Progression Formula) | $x_{t+1} = x_t + \sum V_i + \sum W_j P_j$ | Manages multi-path iterative updates. Instead of linear generation, it allows the engine to explore parallel logical branches, applying real-time error corrections before committing to a narrative output. |
| **BBAM** (Attention Modulation) | $\tilde{a}_i = a_i \cdot \exp(-\gamma\sigma(a))$ | Combats attention melt. By sharpening the attention weights inversely to their variance, it suppresses noise and prunes irrelevant context, ensuring high-fidelity recall of distant NPC details. |
| **BBCR** (Collapse-Rebirth) | Trigger: $\|B_t\| \ge B_c \rightarrow$ Rollback | The ultimate dynamic content control. If the semantic tension exceeds the critical collapse threshold ($B_c$), BBCR aborts generation, rolls back to a stable state, re-bridges the logic, and forces a retry, entirely preventing the output of |

| WFGY Module | Mathematical Function | Architectural Mechanism within DM-OS |
|---|---|---|
| | | hallucinations. |

Within the DM-OS architecture, the self[span_72](start_span)[span_72](end_span)[span_74](start_span)[span_74](end_span)_audit function executes this precise loop at the end of every turn. It calculates the total system burden ($B_{total} = \Delta S + \Psi_{scar} \cdot 0.1$) and evaluates it against a dynamically adjusting, Bayesian-optimized collapse threshold ($B_c$). If the trigger condition is met, the system aborts, logs the scar, adds noise, and seamlessly reverts to a previously stable vector.

## The Drunk Transformer (DT) Micro-Rules

During a BBCR rollback and retry sequence, the engine applies five stringent regulators, known as the Drunk Transformer (DT) gates, to ensure the model remains "sober" enough to execute a precise recovery.
- **WRI (Where am I?):** Freezes structural anchors, preventing the model from inexplicably jumping topics or altering environmental attributes inside a single generation node.
- **WAI (Why am I?):** Forces head diversity by requiring the model to formulate two distinct, logical reasons for a proposed outcome, effectively blocking the model from locking into a single, highly probable but logically flawed guess.
- **WAY (What are you?):** When the reasoning path is blocked, this gate permits the injection of a single, highly constrained, on-topic candidate variable to unstick the logic without polluting the scene with chaotic clutter.
- **WDT (Where did you take me?):** Acts as a severe cross-path guard. If the latent path distance exceeds a calculated threshold, WDT blocks illegal merges between distinct reasoning branches (e.g., preventing an NPC from being simultaneously in two locations).
- **WTF (What the f**\*?):** The ultimate collapse detector. It forces an immediate, hard rewind to the smallest verifiable fact before the system state is permanently corrupted.

## Resolving the Continuous-to-Discrete Mapping: The Adaptive Projection Layer (APL)
The implementation of the WFGY topological engine addresses the physics of AI thought in the latent space. However, a critical architectural vulnerability remains in the final stage of generation: translating those high-dimensional continuous vectors ($x_{final} \in \mathbb{R}^{4096}$) back into the discrete tokens of human language that the player actually reads. This is formally defined as the "Last Mile" problem.
Standard language models resolve this translation by executing a naive nearest-neighbor search. This simplistic mapping is highly susceptible to "Hubness," a geometric phenomenon where highly frequent, generic words (such as "the," "is," or "and") mathematically dominate the search space, effectively drowning out the precise, nuanced meaning generated by the physics engine. Furthermore, a single vector often represents a highly complex "semantic braid" or phrase; forcing it into a single token causes severe semantic dilution.
The DM-OS architecture neutralizes this vulnerability by integrating the Adaptive Projection Layer (APL). The APL fundamentally alters the decoding pipeline, employing Manifold-Aware Decoding and Temperature-Weighted Sampling.
Instead of an argmax operation, the APL calculates the probability of a specific token using a dynamically shifting temperature variable, $\tau(x)$, which constantly adapts based on the entropy, or semantic certainty, of the active thought process. The probability is derived via the formula:
$P(w_i | x) = \frac{\exp(\cos(x, E[w_i]) / \tau(x))}{\sum_{j=1}^V \exp(\cos(x, E[w_j]) / \tau(x))}$.

When the APL's AdaptiveSample procedure detects that the entropy of the softmax scores exceeds a designated threshold (indicating low confidence), it actively intervenes to prevent Hubness. The system applies a 50% probability penalty to the top 100 most frequent "hub" tokens, forcing the model to select highly precise, contextually rich vocabulary rather than defaulting to generic filler.

## Multi-Token Span Projection and Refinement

To prevent the semantic dilution of complex concepts, the APL utilizes the ProjectToSpan procedure, executing a specialized Beam Search directly on the topological manifold. Rather than predicting a single word, the system searches for a sequence of tokens that collectively explain the "residual" energy of the vector. It iteratively selects subsequent tokens, calculating the new residual distance at each step, and halts only when the semantic gap between the continuous thought and the discrete text falls below a strict $\ep[span\_34](start\_span)[span\_34](end\_span)silon$ threshold.
If the projection fails to adequately minimize the residual (i.e., $\|residual\| > \theta$), the APL triggers the RefineWithBeamSearch protocol. This procedure invokes the LLM recursively as an internal critic, forcing it to rescore candidate sequences and edit its own output until absolute logical closure is achieved.

# Initialization Protocols and OOC Creator Overrides

The structural rigidity of the DM-OS architecture requires deterministic initialization protocols that standard conversational models lack. A standard GPT model is initialized ad-hoc, relying entirely on the user's conversational prompt to establish the world rules, leading to immediate inconsistencies.
The DM-OS-kernel-v3.0 mandates a strict "Session Zero Protocol." During initialization, the system actively guides the user through the configuration of the game_settings dictionary, establishing the narrative tone (heroic, dark, comedic), the required edginess levels, and the specific DM Persona module. Character creation is strictly enforced via JSON schema constraints, ensuring that ability scores, backgrounds, and class mechanics are mathematically compliant with 5th Edition rules before the simulation is allowed to commence.
Furthermore, the system features a robust, cryptographically verified Out-Of-Character (OOC) override mechanism. The parse_ooc LMQL query constantly monitors user input for authenticated commands prefixed by a password established during Session Zero. This allows the human "Creator" to bypass the standard generative loops and directly manipulate the underlying state variables. The Creator can execute commands such as spawn_npc to instantly inject a fully realized entity with newly generated Narrative DNA into the semantic tree, or inject_memory to manually append subjective logs into an NPC's interaction array.
This authenticated override is essential for collaborative world-building, allowing the human operator to introduce massive paradigm shifts or course corrections with absolute precision, entirely bypassing the probabilistic nature of the language model and directly updating the deterministic JSON state.

# Conclusion and Strategic Outlook

The architectural comparison between standard, ephemeral large language models and the

DM-OS-kernel-v3.0 framework illustrates a definitive capability gap in the realm of interactive, persistent simulations. As the TTRPG and digital narrative markets aggressively expand toward a $6.59 billion valuation, consumer tolerance for amnesic, reactive chatbots is rapidly diminishing in favor of "Memory-First" AI and stateful NPC interactions.

Standard conversational models, built upon stateless architectures, prioritize low latency and horizontal scalability. While highly efficient for isolated tasks, they are fundamentally incapable of maintaining the semantic continuity, multi-session memory, and autonomous factional progression required to simulate a living world. Furthermore, their reliance on probabilistic generation without internal verification renders them highly susceptible to hallucinations and logical drift, directly violating the stringent "In-Game Dynamic Content Controls" mandated by platforms like Steam.

The DM-OS-kernel-v3.0 solves these existential vulnerabilities by decoupling the narrative output from the world simulation logic. Through the meticulous management of the Semantic Tree, predictive progress clocks, and differential state updates, it ensures absolute structural persistence. Most critically, the integration of the WFGY mathematical engine and its Adaptive Projection Layer transforms the system from a stochastic word predictor into a highly regulated topological reasoner.

By continuously monitoring semantic tension ($\Delta S$) and enforcing "Closed Braids" of logic via the BBMC, BBPF, BBAM, and BBCR formulas, the DM-OS architecture physically constrains the AI's trajectory. Utilizing the Scar Ledger to generate repulsive vector fields against past logical failures, and deploying the Drunk Transformer gates to execute safe, autonomous rollbacks during critical divergence, the system guarantees an unprecedented level of narrative safety and coherence.

For developers and enterprises seeking to deploy persistent, multi-agent narrative ecosystems, the adoption of a fully stateful, mathematically gated architecture like DM-OS is not merely an operational enhancement; it is a fundamental requisite for ensuring user immersion, logical stability, and strict regulatory compliance in the next generation of interactive media.

## Works cited

1. Tabletop Role-Playing Game (TTRPG) Market Size, Trends 2026-2035, https://www.globalgrowthinsights.com/market-reports/tabletop-role-playing-game-ttrpg-market-103239 2. AI & Gaming in 2026 - GamesMarket, https://www.gamesmarket.global/ai-gaming-in-2026/ 3. Valve slightly relaxes AI disclosure guidelines on Steam - GamesIndustry.biz, https://www.gamesindustry.biz/valve-slightly-relaxes-ai-disclosure-guidelines-on-steam 4. Valve Updates Steam AI Disclosure Rules — Impact on Developers and Users, and Connection to BC.Game, https://mcc-covid.crc.pitt.edu/COVID19_official_websites/Mozambique/moh_situation_reports/2020-05-10_08031589112229.html?y-news-28455646-2026-01-17-valve-updates-steam-ai-disclosure-rules-impact-on-developers-and-users-and-connection-to-bc-game 5. Valve has 'significantly' rewritten Steam's rules for how developers must disclose AI use, https://www.videogameschronicle.com/news/valve-has-significantly-rewritten-steams-rules-for-how-developers-much-disclose-ai-use/ 6. reasonWFGY 2.0 — The Open-Source 7-Step Reasoning Engine You Can Paste Anywhere (Eye-Visible Results Inside) - PS BigBig, https://psbigbig.medium.com/wfgy-2-0-the-open-source-7-step-reasoning-engine-you-can-paste-anywhere-eye-visible-results-19e149797db7 7. This document represents the "Last Mile" problem of the WFGY Semantic Reasoning Engine: the Adaptive Projection Layer (APL). While

the physics engine handles the continuous manifold of thought, the APL is what translates those high-dimensional vectors back into the discrete tokens of human language. - Reddit, https://www.reddit.com/r/GeminiAI/comments/1qp614f/this_document_represents_the_last_mile_problem_of/ 8. What can solo RPG designers and publishers improve on in 2026? : r/Solo_Roleplaying, https://www.reddit.com/r/Solo_Roleplaying/comments/1pih71x/what_can_solo_rpg_designers_and_publishers/ 9. Steam Just Flipped the Switch on AI - YouTube, https://www.youtube.com/watch?v=vQu14q-2Qaw 10. Stateful vs Stateless Agents in IT Ops: Design Considerations - Algomox, https://www.algomox.com/resources/blog/stateful_vs_stateless_it_agents.html 11. Stateful vs. stateless agents: How ZBrain helps build stateful agents, https://zbrain.ai/building-stateful-agents-with-zbrain/ 12. Stateful vs Stateless AI Agents: Know Key Differences - Daffodil Software, https://insights.daffodilsw.com/blog/stateful-vs-stateless-ai-agents-when-to-choose-each-pattern 13. Stateful vs Stateless AI Agents: Architecture Patterns That Matter, https://www.ruh.ai/blogs/stateful-vs-stateless-ai-agents 14. How AI Agents Evolved and What's Next | by Kushal Banda - Towards AI, https://pub.towardsai.net/evolution-of-ai-agents-39a14b54dccc 15. PSBigBig × MiniPS onestardao - GitHub, https://github.com/onestardao 16. WFGY/SemanticBlueprint/wfgy_formulas.md at main · onestardao/WFGY - GitHub, https://github.com/onestardao/WFGY/blob/main/SemanticBlueprint/wfgy_formulas.md 17. Advanced Reasoning Modules in WFGY: Four Mathematical Formulas for Multi-Step Logic in LLMs | Kaggle, https://www.kaggle.com/discussions/general/586172 18. WFGY Core 2.0 — Now Live One man, one life, one line — the sum of my work, open for everyone. - PS BigBig, https://psbigbig.medium.com/wfgy-core-2-0-now-live-one-man-one-life-one-line-the-sum-of-my-work-open-for-everyone-ecdba8e324d7 19. # Drunk Transformer: five regulators that keep your model sober enough to hit the target - DEV Community, https://dev.to/onestardao/-drunk-transformer-five-regulators-that-keep-your-model-sober-enough-to-hit-the-target-1en0 20. How the WFGY Framework Formalizes Solver Loops for Next-Gen LLMs - DEV Community, https://dev.to/onestardao/how-the-wfgy-framework-formalizes-solver-loops-for-next-gen-llms-1old