

# A Technical Blueprint for AI-Driven Reverse Engineering of Game Design

## Part I: An Architectural Blueprint for the Game Design Extractor

This report outlines a comprehensive technical architecture for an AI-driven pipeline designed to reverse-engineer game mechanics from visual data. The system's primary function is to analyze gameplay footage at scale, deconstruct the underlying design principles, and output structured, engine-agnostic prompts capable of driving generative AI systems for behavior creation. This document serves as a strategic and implementational guide for the development of this pipeline, focusing on the core analysis engine that translates raw visual input into actionable design specifications.

### Section 1.1: Deconstructing the End-to-End Vision: From Scraped Pixels to Synthesized Behaviors

The strategic vision for this project transcends simple data analysis; it aims to establish a new paradigm in game development, shifting from purely manual design to AI-assisted reverse-engineering and procedural generation. This process is conceptualized as a four-phase pipeline, where each stage serves a distinct purpose and provides a critical handoff to the next.

The end-to-end workflow is defined as follows:

1. **Phase 1: Game Scraping & Raw Data Extraction:** This initial phase involves the systematic collection of visual data from tens of thousands of games. The existing roam-game-analysis repository serves as the foundation for this effort, providing the tooling to capture screenshots, video clips, and associated metadata [User Query]. The output of this phase is a curated and structured dataset of visual game information.

2. **Phase 2: Game Analysis & Data Application:** This is the core intellectual engine of the pipeline. It ingests the raw visual data from Phase 1 and performs a deep analysis to extract core game design elements. Its primary objective is to deconstruct the "what" and "why" of gameplay interactions and produce a detailed, engine-agnostic JSON output.<sup>1</sup> This report focuses predominantly on the architecture of this phase.
3. **Phase 3: Behavior Trees & Unity Integration:** This phase consumes the JSON specification generated by Phase 2. The JSON acts as a high-level prompt or "source code" for a sophisticated system within the Unity engine, which is responsible for synthesizing complete, functional behaviors, such as behavior trees or state machines [User Query]. The success of this phase is directly dependent on the clarity, completeness, and structural integrity of the JSON output.
4. **Phase 4: Backend Integration & Orchestration:** The final phase involves integrating the entire pipeline into a broader production environment. This includes storing scraped data and generated JSON in queryable databases, creating services to search and retrieve game analyses, and connecting the output to an orchestrator agent responsible for full-scale game and asset generation [User Query].

The interdependencies between these phases are critical. The quality of the scraped data from Phase 1—including standardized clip lengths, resolution, and metadata—directly impacts the efficacy of the analysis in Phase 2. Similarly, the JSON schema produced by Phase 2 must be meticulously co-designed with the requirements of the Unity behavior generation system in Phase 3. This JSON is not merely data; it is a formal specification for generating game logic.<sup>1</sup>

Strategically, this pipeline addresses a fundamental challenge in Procedural Content Generation (PCG). While traditional PCG excels at creating vast amounts of content, it often struggles to produce content that is logically coherent, stylistically consistent, or embodies a specific design intent.<sup>2</sup> This system moves beyond generating random variations by learning the implicit

*rules* of game design from a massive corpus of existing, successful games, enabling the generation of content that is not just novel but also well-designed.<sup>4</sup>

## Section 1.2: The Phase 2 Core: A Multi-Layered Analysis Engine

To effectively deconstruct game design from visual data, the architecture of the Phase 2 analysis engine is best conceptualized as a multi-layered cognitive model. This structure mirrors the process a human expert would undertake to analyze gameplay, providing a robust separation of concerns and enhancing the system's modularity and extensibility. The architecture is composed of three distinct layers:

1. **Layer 1: Perception (The "What"):** This foundational layer is exclusively concerned with visual processing. Its function is to answer the question, "What is happening on screen?" It achieves this by identifying and classifying game objects, tracking their movement over time, reading text and values from the UI, and detecting discrete, atomic gameplay events. This layer directly implements the functionality described as the "Visual Parsing Module" in the project's foundational design document.<sup>1</sup>
2. **Layer 2: Inference (The "Why"):** This layer serves as the reasoning engine of the pipeline. It receives the structured stream of perceived states and events from Layer 1 and seeks to answer the question, "Why did that happen?" This involves deducing the underlying rules, mechanics, and causal relationships that govern the game world. For example, it infers that a player's attack *causes* an enemy's health to decrease, or that an enemy begins chasing the player *because* the player entered a specific range. This layer is the heart of the "LLM Interpretation & Abstraction Module".<sup>1</sup>
3. **Layer 3: Abstraction & Generation (The "How"):** The final layer is responsible for translating the inferred logic and quantified parameters from Layer 2 into the final, structured JSON output. It answers the question, "How can this knowledge be represented for a machine?" This involves generating rich, descriptive text prompts for assets and behaviors, and meticulously populating the predefined JSON schema that will be consumed by the downstream Unity engine.

This three-layer architecture is more than an engineering convenience; it represents a more robust and scalable approach than a monolithic, end-to-end model. A naive approach might involve feeding raw video into a single large model and hoping for a correct JSON output. However, research into complex AI tasks like game playing and program synthesis highlights the limitations of such methods, which often suffer from logical inconsistencies and "hallucinations".<sup>5</sup> The proposed modular structure, which formalizes the separation suggested in the initial project brief <sup>1</sup>, allows for specialized components at each layer.

The Perception layer can utilize highly optimized computer vision models tailored for

tasks like object detection and tracking. The Inference layer can employ a different set of tools, such as probabilistic models, logic solvers, or LLMs fine-tuned specifically for causal reasoning. The Abstraction layer can then use another class of models focused on structured data generation and high-quality natural language. This modularity ensures that each component can be developed, tested, and upgraded independently. If a superior object detection model becomes available, only the Perception layer needs to be updated. If the quality of the natural language descriptions needs improvement, the LLM in the Abstraction layer can be retrained without altering the core inference logic. This design philosophy makes the entire system more maintainable, debuggable, and prepared for future advancements.

## **Part II: The Perception Layer: Vision-Based Game State Deconstruction**

The Perception Layer is the sensory input for the entire analysis pipeline. Its mission is to transform the raw, unstructured pixel data from screenshots and video clips into a rich, structured, time-series representation of the game state. This involves a suite of computer vision tasks working in concert to identify what entities are present, how they are moving, what the explicit state of the game is, and what discrete actions are taking place.

### **Section 2.1: Entity and Object Recognition: Identifying the Actors**

The first step in understanding a game scene is to identify all the relevant entities within it. This includes player characters, various types of enemies, non-player characters (NPCs), projectiles, collectible items, and interactive environmental objects.

**Methodology:** This task will be accomplished using state-of-the-art object detection models. The recommended approach is to leverage the YOLO (You Only Look Once) family of models, such as YOLOv11 or its successors. These models are renowned for their exceptional balance of high speed and strong accuracy, making them ideal for processing video streams where real-time or near-real-time performance is a

consideration.<sup>7</sup> Their architecture is well-suited to detecting objects of varying sizes and aspect ratios, a common characteristic of game scenes with diverse characters and items.<sup>8</sup>

The implementation requires training a custom object detection model on a large, annotated dataset of game screenshots. The roam-game-analysis repository can be leveraged to build the necessary infrastructure for data collection and annotation. For each detected object, the model will output a structured data packet containing:

- A temporary unique ID for tracking purposes.
- A class label (e.g., PlayerCharacter, Enemy\_Ranged, Collectible\_Coin).
- A confidence score for the detection.
- A bounding box in normalized screen coordinates [x, y, width, height].

## Section 2.2: Spatiotemporal Analysis: Tracking Movement and Kinematics

Static object detection is insufficient; understanding behavior requires tracking entities across time to analyze their movement. This spatiotemporal analysis provides the raw data for inferring kinematic properties and movement patterns.

**Methodology:** This will be implemented using a "tracking-by-detection" framework, which is a standard and effective approach in multi-object tracking (MOT).<sup>9</sup> In this paradigm, the object detector from Section 2.1 is run on each frame of a video clip. Subsequently, a tracking algorithm associates these detections across frames to form persistent "tracklets" for each unique entity.

For this association task, algorithms like DeepSORT, which combines motion prediction with appearance features, or simpler, faster methods like a Kalman Filter, can be employed.<sup>9</sup> A Kalman Filter is particularly effective for predicting an object's position in the next frame based on its current state (position and velocity), making the association process more efficient.

Once these tracklets are established, a wealth of kinematic data can be derived for each entity:

- **Position:** The centroid of the bounding box at each time step.
- **Velocity:** The vector change in position between frames.
- **Acceleration:** The rate of change of velocity.

- **Trajectory:** The sequence of positions over a given time window.

For more granular analysis of complex animations or non-rigid objects, optical flow techniques can be used to estimate the motion vector of every pixel or pixel block between frames, providing a dense motion field of the entire scene.<sup>9</sup>

The output of this module is a time-series of kinematic data for each tracked entity, which will directly populate the `estimated_transform` and inform the `inferred_kinematics` fields in the final JSON output.

### Section 2.3: HUD and UI Analysis: Quantifying the Game State

Many critical game mechanics, such as health, damage, ammunition, and score, are explicitly represented in the game's Heads-Up Display (HUD). Extracting this information is paramount for quantifying game rules. This process can be achieved without access to the game's source code by analyzing the visual data stream.<sup>12</sup>

**Methodology:** The approach will adapt and modernize the robust, code-free techniques detailed in existing research on visual game metrics analysis.<sup>12</sup> The process involves several stages:

1. **HUD Localization:** First, a classifier is trained to identify the static regions of the screen that contain the HUD. This can be done with a simple Convolutional Neural Network (CNN) or a cascade classifier as described in the source material. This initial step makes all subsequent analysis robust to changes in HUD position, scale, or aspect ratio across different games or user settings.
2. **Element Segmentation:** Within the localized HUD region, more specific elements like health bars, ammo counters, and mini-maps are identified. This can be done using template matching with pre-defined icons or by training smaller, specialized object detectors for these UI components.
3. **Value Extraction from Bars:** For progress bars (representing health, mana, stamina, etc.), color thresholding within the HSV (Hue, Saturation, Value) color space is a highly effective technique. By isolating pixels of a specific color (e.g., bright red for health), the system can create a mask for the "filled" portion of the bar. The value can then be calculated as the ratio of the filled area's width or height to the total bar container's width or height. This directly enables the inference of parameters like `max_health_inferred_from_ui` and allows for tracking observed\_changes like a percentage decrease in health.<sup>1</sup>

4. **Value Extraction from Text:** For numerical displays like ammo counts, scores, or timers, an Optical Character Recognition (OCR) engine is employed. While Tesseract is a viable option as demonstrated in the research <sup>12</sup>, more modern deep learning-based OCR systems may offer higher accuracy. Crucial pre-processing steps include isolating the text region, converting it to grayscale, and applying a binary threshold to create a clean, high-contrast image of the characters before feeding it to the OCR engine.

The output of this module is a structured, time-stamped stream of UI data, such as { "ui\_element\_id": "player\_health\_bar", "timestamp": 12.1, "value\_percentage": 75.0 }.

## Section 2.4: Event Recognition: Identifying Atomic Gameplay Actions

To understand game mechanics, the system must recognize not just objects but also the actions they perform. This involves detecting and labeling discrete, meaningful gameplay events as they occur in the video, such as `player_attacks`, `enemy_takes_damage`, `item_collected`, or `player_jumps`.

**Methodology:** This task requires models capable of understanding temporal dynamics in video, as an event is defined by changes over time.

- **Two-Stream Convolutional Networks:** Architectures like I3D (Inflated 3D ConvNets) are well-suited for this task. They process two parallel input streams: one for RGB frames, capturing spatial appearance (the look of the character and weapon), and one for optical flow, capturing temporal motion (the swinging of a sword or recoil of a gun). By fusing the features from both streams, the model can make highly accurate action classifications.<sup>11</sup>
- **Video Transformers:** More recent architectures like VideoMAE or the multi-modal X-CLIP leverage the power of transformers for video understanding.<sup>14</sup> These models can capture complex, long-range spatiotemporal relationships and have demonstrated state-of-the-art performance on action recognition benchmarks. Given their power, they are a primary candidate for this module.

Significant data collection and annotation are required to train these models. This involves creating a dataset of thousands of short, labeled video clips, with each clip corresponding to a specific game event.

The output of this module is a time-stamped log of detected events, structured to



include the event type and the entities involved, for example: (timestamp: 15.4, event\_type: 'ENEMY\_TAKES\_DAMAGE', source\_id: 'player\_character\_01', target\_id: 'enemy\_melee\_grunt\_003'). This event stream is a primary input for the Inference Engine.

The individual modules of the Perception Layer should not be viewed as independent, parallel processes. Instead, they form a symbiotic system where the output of one module can critically inform and enhance the performance of another. A truly advanced pipeline will be architected as a directed graph of information flow, not as a set of isolated tasks. For instance, knowing which objects are present in a scene (from object detection) dramatically constrains the set of possible events that the event recognition model needs to consider. If no "projectile" object is detected, the likelihood of a "projectile\_fired" event is virtually zero. This contextual information can be fed as an additional input to the event recognition model, improving its accuracy and reducing false positives.

Furthermore, a detected event can trigger a focused analysis in another module. The detection of an ENEMY\_ATTACK\_HIT\_PLAYER event should be immediately followed by a change in the player's health bar value from the HUD analysis module. By correlating these two data streams, the system can validate its own perceptions. If a health drop is observed without a corresponding damage event, it may signal the presence of a previously unknown mechanic, such as poison or environmental damage, providing a rich signal for the Inference Engine. This cross-validation creates a self-correcting loop that improves the overall fidelity of the game state representation. The final output of the Perception Layer should therefore be a single, unified "game state frame" at each time step, integrating object lists, their kinematics, UI values, and active events into one cohesive data structure. This rich, contextualized data stream is exponentially more valuable to the Inference Engine than four separate, uncorrelated streams would be.

## **Part III: The Inference Engine: From Observation to Causal Mechanics**

This part of the pipeline represents the most significant research and development challenge. It is the cognitive core of the system, designed to move beyond simple observation to genuine understanding. Its function is to take the structured stream of



"what" is happening from the Perception Layer and synthesize the "why"—the underlying game mechanics, rules, and causal relationships that govern the gameplay.

### Section 3.1: A Framework for Causal Inference in Gameplay

The primary objective of the Inference Engine is to establish causal links between events and changes in the game state. It must distinguish true causation from mere correlation. For instance, an enemy's health bar decreasing *at the same time* as a player's attack animation is a correlation. The goal is to infer that the attack *caused* the health decrease. Observational data from gameplay is inherently noisy and full of confounding variables, making this a non-trivial task.<sup>15</sup>

**Methodology:** To tackle this, the engine will employ a "Hypothesize-and-Test" framework grounded in the principles of causal inference from observational data.<sup>15</sup>

1. **Event-State Change Pairing:** The engine continuously monitors the data stream from the Perception Layer, looking for temporal correlations. It logs pairs of (Event, State Change) that occur within a very short time window (e.g., a few frames). An example pair would be (Event: 'player\_attack\_animation\_ends', StateChange: 'enemy\_health\_decreased').
2. **Hypothesis Generation:** For every pair that occurs with a frequency above a certain threshold, the system generates a formal causal hypothesis. This hypothesis takes the form: Event(E) with parameters (P\_E) CAUSES StateChange(S) with parameters (P\_S). For the example above, the hypothesis would be: player\_attack(source: player, target: enemy) CAUSES enemy\_health\_change(amount:  $\Delta$ ).
3. **Counterfactual and Confounder Analysis:** This is the critical step. The system actively searches the entire dataset for evidence that could support or refute the hypothesis. This is analogous to searching for natural control groups in an observational study.<sup>15</sup>
  - **Supporting Evidence:** The system finds all other instances where the player attacked that specific type of enemy, and the enemy's health decreased.
  - **Refuting Evidence (Counterfactuals):** The system looks for cases where the player attacked, but the enemy's health did *not* decrease. It then analyzes the conditions of these cases. For example, it might discover that in all such instances, the distance between the player and enemy was greater than 2 units. This leads to a refinement of the hypothesis: player\_attack(...) CAUSES...

IF distance\_to\_target < 2.

4. **Parameter Quantification:** Once a causal link is established with high confidence, the system aggregates the data from all validated instances to quantify its parameters. By analyzing the distribution of health decreases ( $\Delta$ ) across all successful attacks, it can infer the inferred\_player\_attack\_damage. For example, if the health drops are consistently ``, the system can confidently infer the damage value is 20.<sup>1</sup> Similarly, if the event collectible\_coin\_appears consistently follows the event enemy\_X\_is\_destroyed, the system infers the spawn rule coin\_on\_enemy\_X\_defeat.<sup>1</sup>

### Section 3.2: Rule Induction for Behavioral Logic: Synthesizing "IF-THEN" Statements

Beyond quantifying simple parameters, the Inference Engine must deconstruct the complex logic that drives entity behavior. This task can be framed as a form of program synthesis: the system observes an entity's actions and synthesizes the "program" (the set of rules) that governs its AI.<sup>18</sup>

**Methodology:** The approach involves inferring a Finite State Machine (FSM) or a set of behavior tree rules from the observed action sequences.<sup>20</sup>

1. **State Identification:** The engine first needs to identify the distinct behavioral "states" an entity can occupy. These states are not explicit but can be clustered and inferred from the entity's kinematic data (from Perception Layer Section 2.2) and the events it participates in (Section 2.4). For example, periods of slow, back-and-forth movement can be clustered into a Patrolling state. Periods of high-velocity movement directly toward the player character can be clustered into a Chasing state. Periods of stationary behavior punctuated by attack events can be clustered into an Attacking state.
2. **Transition Trigger Discovery:** The system then focuses on the moments of transition between these inferred states (e.g., the exact frame an entity stops Patrolling and starts Chasing). It analyzes the complete game state at that moment to identify the conditions that were met immediately prior to the transition. These conditions form the "IF" clause of a behavioral rule.
3. **Rule Synthesis Example:**
  - **Observation:** The system identifies a transition from an inferred Patrolling state to a Chasing state for enemy\_grunt\_003.

- **Condition Analysis:** It examines the game state in the frames just before the transition and finds that the player's position moved to within a 5-unit radius of the enemy. In other instances of this same transition, the same condition held true.
- **Rule Induction:** The system synthesizes the rule: IF current\_state IS 'Patrolling' AND distance\_to\_player < 5.0 THEN transition\_to\_state('Chasing').
- This process is repeated for all observed state transitions (e.g., Chasing to Attacking, Chasing to Patrolling if the player moves out of range), gradually building a complete and robust FSM that describes the entity's logic.

### Section 3.3: Inferring Latent and Implicit Parameters

Many crucial game mechanics are not directly visible in the UI and must be inferred from the timing and frequency of events. These "latent" parameters include things like attack cooldowns, weapon reload times, ability charges, or enemy spawn rates.

**Methodology:** This relies on statistical analysis of the time-stamped event stream generated by the Perception Layer.

- **Time Interval Analysis:** To infer an attack cooldown, the system will query all attack events initiated by a specific entity type. It will then calculate the time delta,  $\Delta t$ , between each consecutive pair of attacks. By plotting a histogram of these  $\Delta t$  values, a clear peak will emerge around the true cooldown time (e.g., a cluster of values around 2.0 seconds for an attack with a 2-second cooldown). The mean or median of this cluster becomes the inferred\_attack\_cooldown.
- **Probabilistic and Unsupervised Methods:** To infer mechanics like spawn rates, the system can use unsupervised learning techniques.<sup>21</sup> By observing a specific region of a level over many gameplay sessions, it can count the number of enemies that appear and the intervals between their appearances. This data can be used to fit a probability distribution, inferring parameters like spawn\_rate\_per\_minute or avg\_spawn\_count.
- **Adapting Game Theory Concepts:** For more complex, strategic interactions, concepts from game theory and reinforcement learning can be adapted. By treating observed gameplay segments as "playouts" in a Monte Carlo simulation, the system can analyze the outcomes of different action sequences. This can help infer the game's implicit reward structure, revealing which actions or states are considered "good" by the game's design, even without an explicit score.<sup>22</sup>

A critical distinction of this entire inference approach is its fundamental difference from behavioral cloning (BC). BC, a form of imitation learning, aims to learn a direct mapping from observation to action, essentially creating a policy  $\pi(\text{action} \mid \text{observation})$  that mimics a human player.<sup>23</sup> However, BC is notoriously brittle and susceptible to "causal confusion"—learning spurious correlations instead of true causal links.<sup>23</sup> For example, a BC agent might learn to fire its weapon when it sees a muzzle flash (the

*effect* of firing) rather than when it sees an enemy (the *cause* for firing). This is because it does not model the underlying dynamics of the world.

The Inference Engine described here is explicitly designed to avoid this pitfall. Its goal is not to clone the player's policy, but to reverse-engineer the *game's mechanics*. It uses the player's actions as a series of experiments or probes to discover the rules of the environment. The system models the world's transition function,  $P(\text{next\_state} \mid \text{state}, \text{action})$ , which is inherently more robust and closer to the ground-truth game logic. When the player performs an "attack" action, the system doesn't just learn to associate the current screen with that action; it observes the *consequences* of that action on the entire game state—changes in enemy health, score, inventory, etc. By focusing on these state transitions, the system infers the game's underlying rules, creating a model of the game itself, not just a shallow copy of one player's behavior. This makes the resulting knowledge far more generalizable and useful for generating new, coherent game content.

## **Part IV: The Abstraction & Generation Layer: Structuring Knowledge for Generative AI**

The final layer of the Phase 2 pipeline is responsible for communication. It takes the complex, inferred knowledge from the Inference Engine—the causal links, behavioral rules, and quantified parameters—and translates it into a clear, structured, and actionable format. This involves populating the definitive JSON schema and generating high-quality natural language prompts that are optimized for consumption by downstream generative AI systems in Unity.

Section 4.1: The Unified JSON Schema: A Deep Dive

The JSON output is the primary deliverable of the analysis pipeline and serves as the formal contract between the analysis and generation systems. A robust, comprehensive, and well-documented schema is therefore essential for the project's success. The following schema builds upon the illustrative example provided in the initial project documentation <sup>1</sup>, expanding it with fields derived from the advanced capabilities of the Perception and Inference layers.

A key enhancement is the evolution of the behavior description. Instead of a single text prompt, the schema now includes a `behavioral_logic` object. This object provides a structured, machine-readable representation of the inferred FSM or rule set. The `behavior_description_prompt` is then generated as a natural language *summary* of this structured logic, providing the best of both worlds: precise data for programmatic interpretation and a rich, contextual prompt for LLM-driven generation.

Table: Detailed Unified JSON Schema Definition

Key Path	Data Type	Description & Rationale	Example
game_source_id	String	Unique identifier for the source game/clip being analyzed.	"archero_clip_01"
analysis_timestamp	String (ISO 8601)	Timestamp of when the analysis was performed.	"2024-10-26T14:30:00Z"
global_inferences	Object	Scene-wide observations.	
global_inferences.inferred_genre_tags	Array of Strings	Inferred genre tags based on observed mechanics.	``
global_inferences.camera_perspective_approx	String	Estimated camera perspective.	"TopDown_SlightlyAngled"
scene_elements	Array of Objects	An array of all	

		significant entities identified in the scene.	
scene_elements.element_id	String	A persistent, unique identifier for the entity instance.	"enemy_melee_grunt_003"
scene_elements.element_type	String (Enum)	The high-level type of the entity.	"Enemy"
scene_elements.subtype_inferred	String	A more specific, inferred subtype for the entity.	"MeleeGrunt"
scene_elements.visual_description_prompt	String	A rich text prompt for an image/3D generation model.	"A bulky, red-skinned demonic minion with two large pincers for hands. Slow but determined movement. Low-poly, stylized art style."
scene_elements.estimated_transform	Object	Approximate spatial data at a representative moment.	
scene_elements.estimated_transform.position_2d_normalized	Array of Numbers	Normalized [x, y] screen coordinates.	[0.7, 0.6]
scene_elements.estimated_transform.scale_approx	Array of Numbers	Approximate normalized [width, height] scale.	[0.08, 0.08]
scene_elements.behavioral_logic	Object	<b>Structured, machine-readable behavioral logic.</b>	
scene_elements.behavioral_logic.states	Array of Objects	The inferred states of the entity's FSM.	``
scene_elements.behavioral_logic.transitions	Array of Objects	The rules governing transitions between states.	[{"from": "Patrol", "to": "Chase", "conditions": [{"type": "distance", "target":

			"Player", "op": "<", "value": 5.0}}}]
scene_elements.behavior_description_prompt	String	<b>Natural language summary of the behavioral logic.</b>	"This enemy patrols slowly between two points. If a player enters its visual range (approx. 5 units), it will charge towards the player. Once in melee range (approx. 1 unit), it attacks."
scene_elements.game_mechanic_params	Object	Quantifiable game design values inferred for this entity.	
scene_elements.game_mechanic_params.health	Number	Inferred health points, often from UI analysis.	100
scene_elements.game_mechanic_params.inferred_attack_damage	Number	Inferred damage per attack.	15
scene_elements.game_mechanic_params.inferred_attack_cooldown_sec	Number	Inferred time in seconds between attacks.	2.5
scene_elements.game_mechanic_params.inferred_movement_speed_chase	Number	Inferred speed in units/sec during chase state.	2.0
scene_elements.on_defeat_spawns_inferred	String	A description of what happens when this entity is defeated.	"Upon defeat, this enemy type typically drops 2 to 4 gold coins."
ui_elements_observed	Array of Objects	Data extracted from observed HUD/UI elements.	
ui_elements_observe	String	Unique ID for the UI	"player_health_bar_0



d.ui_element_id		element.	1"
ui_elements_observe d.observed_changes	Array of Objects	A log of observed changes to this UI element's value.	``

## Section 4.2: Generating Visual Asset Prompts

A key output of the system is the ability to generate descriptive text prompts that can guide generative models (like text-to-image or text-to-3D tools) in creating new game assets.<sup>25</sup> This is an "image-to-text" or visual captioning task.

**Methodology:** The cropped image of a detected entity from the Perception Layer serves as the primary input. This image is fed into a multi-modal model, such as one based on the CLIP architecture or a dedicated visual transformer (ViT) paired with a text decoder like GPT-2. The model's task is to generate a textual description of the entity's appearance.

This process is significantly enhanced by injecting contextual information from the Inference Engine. The raw caption generated from the image can be programmatically prepended with inferred tags and details. For example, if the Inference Engine has classified an entity with subtype\_inferred: "MeleeGrunt" and the game's style is identified as art\_style: "PixelArt", the final prompt can be constructed as: "Pixel art sprite of a Melee Grunt enemy. It is a small, round, red demonic creature with short horns and glowing yellow eyes." This fusion of direct visual description with inferred game-specific context results in a far more effective and targeted prompt for asset generation.<sup>28</sup>

## Section 4.3: Natural Language Generation for Behavioral Prompts

The behavior\_description\_prompt is a critical piece of the output, as it must fluently and accurately communicate the complex, inferred behavioral logic to a downstream AI. This is a classic Natural Language Generation (NLG) task, specifically of the "data-to-text" variety, where structured data is converted into a human-readable narrative.<sup>29</sup>

**Methodology:** The primary input for this module is the structured behavioral\_logic object from the JSON schema, which contains the inferred FSM, transition rules, and associated parameters like speed and range.

- **LLM-based Generation:** The most flexible and powerful approach is to use a fine-tuned Large Language Model (LLM), such as T5 or a model from the GPT family.<sup>31</sup> This model would be trained on a dataset of paired examples, where each example consists of a structured logic JSON and its corresponding high-quality, human-written description. The model learns to "translate" the formal logic into a natural narrative.
- **Template-based Fallback:** For simpler behaviors or as a highly reliable baseline, a template-based system can be implemented. This system would use predefined sentence structures and fill in the blanks with the inferred parameters. For example: "This entity primarily follows a {state} behavior. When {condition}, it transitions to the {next\_state} state, moving at a speed of {speed} units per second."

### Example Translation:

- **Input Structured Logic:**

```
JSON
{
  "states": [{"name": "Patrol"}, {"name": "Chase"}, {"name": "Attack"}],
  "transitions": [
    {"from": "Patrol", "to": "Chase", "conditions": [{"type": "distance_to_player", "op": "<", "value": 5.0}]},
    {"from": "Chase", "to": "Attack", "conditions": [{"type": "distance_to_player", "op": "<", "value": 1.0}]}
  ],
  "params": {"patrol_speed": 1.0, "chase_speed": 2.0, "attack_damage": 15}
}
```

- **Generated Natural Language Prompt:**  
"This enemy patrols slowly at approximately 1.0 units/sec. If a player character enters its visual range (approx. 5 units), it will stop patrolling and charge directly towards the player at a speed of 2.0 units/sec. Once within close melee range (approx. 1 unit), it performs a melee attack that deals 15 damage." This output mirrors the quality and detail level specified in the project's requirements.<sup>1</sup>

Section 4.4: Advanced Prompt Engineering Strategies

The effectiveness of a prompt is not just in its grammatical correctness, but in its ability to elicit the desired, high-quality output from the target AI model.<sup>32</sup> Therefore, the generation of the

behavior\_description\_prompt must be treated as a deliberate engineering task, incorporating proven strategies to maximize its impact on the Unity-side generative system.<sup>33</sup>

**Methodology:** The NLG module from Section 4.3 will be designed to structure its output using the following prompt engineering principles:

- **Chain-of-Thought (CoT) Structuring:** Instead of providing a single monolithic block of text, the prompt can be structured to guide the downstream model through a logical sequence of steps. By breaking down the behavior into its constituent parts, the prompt encourages more structured reasoning from the generation model.<sup>34</sup>
- **Role-Playing:** The prompt can begin by assigning a specific role to the AI, which primes it to access the relevant parts of its training. For example, starting the prompt with "You are an expert game AI designer creating a behavior tree..." sets a clear context.
- **Explicit Parameterization:** All quantified data extracted by the Inference Engine should be explicitly listed within the prompt. This grounds the generative model in concrete facts and prevents it from hallucinating or using arbitrary values.
- **Context and Constraints:** The prompt should clearly define the context, constraints, and goals of the behavior.

The following table provides a practical guide for implementing these techniques, contrasting basic prompts with engineered prompts to illustrate the increase in clarity and effectiveness.

Table: Prompt Engineering Techniques for Behavior Generation

Technique	Basic Prompt Example	Engineered Prompt Example (with rationale)
Role-Playing	"The enemy chases the player and attacks."	"As a behavior tree architect, your task is to

		<b>design the logic for a 'Melee Grunt' enemy.</b> " (Rationale: Assigns a specific expert persona to the LLM, focusing its response.)
<b>Chain-of-Thought</b>	"It patrols, then chases if the player is close, then attacks."	<b>"The behavior is composed of three primary states: 1. PATROL: The enemy moves between predefined points. 2. CHASE: The enemy moves towards the player. 3. ATTACK: The enemy performs a melee strike."</b> (Rationale: Breaks the problem down into logical steps, making the logic easier for the model to follow and implement correctly.)
<b>Explicit Parameterization</b>	"It moves fast when chasing and its attack is strong."	<b>"Use the following exact parameters for the implementation: patrol_speed: 1.0, chase_speed: 2.5, detection_range: 8.0, attack_range: 1.5, attack_damage: 20, attack_cooldown: 2.0 seconds."</b> (Rationale: Provides concrete, non-negotiable data points, preventing the model from inventing its own values and ensuring the generated behavior matches the analysis.)
<b>Context &amp; Constraints</b>	"Make the enemy flee if it gets hurt."	<b>"The enemy has a 'Flee' behavior that triggers under specific conditions. Constraint: This Flee state should only be entered if the enemy's health drops below 25% AND it is outnumbered by more than one player-aligned</b>

		<b>character. It should not flee otherwise."</b> (Rationale: Provides not just the trigger, but also the negative constraints, leading to more nuanced and accurate logic.)
--	--	---

By systematically applying these techniques, the Abstraction & Generation Layer produces prompts that are not just descriptive but are precision instruments designed to control and guide the creative output of the final behavior synthesis engine.

## Part V: Ensuring Robustness, Consistency, and Scalability

For the game design extraction pipeline to be a viable production tool, it must be more than just clever; it must be robust, reliable, and capable of operating at a massive scale. This final section addresses the critical system-wide challenges of maintaining logical consistency, implementing rigorous validation, and architecting for the analysis of tens of thousands of diverse games.

### Section 5.1: The Consistency Challenge: A "Model as a Game" Approach

A significant failure mode for purely generative models is their tendency to produce outputs that violate fundamental game rules, breaking numerical and spatial consistency.<sup>6</sup> A model might generate a score of "100" in one frame and "500" in the next for no reason, or it might render a completely different room when a player turns around and looks back. Such inconsistencies shatter the illusion of a coherent game world.

**Methodology:** To combat this, the pipeline architecture will be augmented with principles from the "Model as a Game" (MaaG) paradigm.<sup>6</sup> This approach involves externalizing critical state information into explicit, non-generative modules that serve as a "ground truth." The generative components of the pipeline are then conditioned on this ground truth, forcing them to remain consistent.

- **Numerical Consistency Module:** An external state tracker will be implemented

to maintain the canonical values of all quantified game variables (e.g., player health, score, ammo count, resources). When the Inference Engine identifies a causal event, such as `player_collects_coin`, it does not try to make a generative model guess the new score. Instead, it issues a deterministic command to this external module, such as `increment_score(value: 1)`. Subsequently, when the Abstraction Layer generates a visual description or a new image frame, it will query this trusted module for the current score and use that value as a hard constraint on the generation process. This architecture, mirroring the LogicNet and External Numerical Record concepts, decouples game logic calculation from visual rendering, ensuring numerical integrity.<sup>36</sup>

- **Spatial Consistency Module:** For games involving exploration across a persistent world, a similar external module will maintain a canonical map of all explored areas. As the Perception Layer observes new locations, they are added to this global map. When the system needs to generate data or descriptions for a previously visited area, it retrieves the authoritative data from this map. This prevents spatial discontinuities, such as objects disappearing or the layout of a room changing upon re-entry. This method directly implements the External Map and Sliding Window Matching techniques, which use pattern matching (like Peak Signal-to-Noise Ratio) to align newly generated frames with the existing map, ensuring seamless spatial coherence.<sup>36</sup>

## Section 5.2: A Framework for Validation and Quality Assurance

To trust the output of the pipeline, a multi-faceted validation and quality assurance framework is essential. This framework will test the system's accuracy, logical coherence, and plausibility.

### Methodology:

- **Ground Truth Comparison:** For a curated set of "gold standard" games—either simple, custom-built test cases or open-source titles where the source code is available—the system's output can be directly validated. The inferred parameters from the JSON output (e.g., `inferred_attack_damage: 15`) can be compared against the actual hardcoded values in the game's source code. This provides a quantitative metric for parameter extraction accuracy.
- **Automated Logical Coherence Testing:** The generated JSON files can be subjected to a suite of automated logical tests. These scripts would check for

internal contradictions, such as a behavior that defines a transition from a state back to itself with no triggering condition, a negative value for a health parameter, or a behavioral\_logic FSM with unreachable states. These tests act as a linter for game design logic.

- **Human-in-the-Loop (HITL) Evaluation:** Ultimately, game design has subjective components. The final validation step involves expert human review. A panel of experienced game designers will be presented with the generated JSON outputs, particularly the behavior\_description\_prompt and the structured behavioral\_logic. They will rate the outputs on criteria such as plausibility, completeness, accuracy, and nuance. This qualitative feedback is invaluable for refining the inference and generation models, drawing parallels to studies where expert gamers were used to quantify player expertise more accurately than self-reporting.<sup>37</sup>

### Section 5.3: Strategies for Scalable Game Analysis

The project's ambition to analyze "10,000s of games" necessitates an architecture designed for scalability, capable of handling immense diversity in genre, visual style, and data volume.

**Methodology:** A multi-pronged strategy is required to achieve this scale:

- **Genre-Specific Analysis Pipelines:** A single, monolithic model attempting to understand every game genre from 2D platformers to 4X strategy games is unlikely to succeed. A more effective approach is to train a fleet of specialized models. A high-level classifier will first categorize a new game by genre (e.g., "First-Person Shooter," "2D Platformer," "Top-Down RPG"). The game's visual data will then be routed to a dedicated analysis pipeline whose Perception and Inference models have been specifically trained on data from that genre. The platformer model will excel at recognizing jumps, moving platforms, and collectibles, while the FPS model will be an expert in analyzing HUDs, weapon switching, and aiming mechanics.
- **Structured Data Management and Querying:** As specified in the project's future goals, all scraped visual data and the resulting JSON analysis outputs must be stored in a centralized, structured database system. This database should be indexed and queryable by game title, genre, inferred mechanics (e.g., "all games with a 'parry' mechanic"), specific parameter values, and other metadata. This will transform the output from a collection of individual files into a powerful,



large-scale analytical resource for discovering design trends.

- **Distributed Cloud-Based Processing:** The analysis of each game is an independent, embarrassingly parallel task. The entire pipeline should be architected to run on a distributed computing framework like Kubernetes, deployed on a cloud platform (e.g., AWS, GCP, Azure). This will allow for the dynamic scaling of GPU-equipped worker nodes, enabling the parallel processing of hundreds or thousands of games simultaneously to meet the throughput requirements.
- **Unsupervised Pre-training for Generalization:** To handle the vast visual diversity of games (from pixel art to photorealism), the vision models in the Perception Layer should undergo extensive pre-training on a massive, unlabeled dataset of general game imagery. Using self-supervised learning techniques, the models can learn a robust internal representation of common visual concepts in gaming before being fine-tuned on specific, labeled tasks like object detection or event recognition.<sup>21</sup> This pre-training step dramatically improves a model's ability to generalize to new, unseen games and art styles, reducing the amount of labeled data needed for each new genre.

## Conclusion and Recommendations

The proposed AI pipeline represents a significant technological undertaking with the potential to fundamentally alter the landscape of game design and development. By systematically deconstructing existing games into their core mechanical and aesthetic components, this system can create an unprecedented, queryable library of game design knowledge. This knowledge, in turn, can fuel generative AI systems to create novel, coherent, and well-designed behaviors and content at scale.

The architectural philosophy outlined in this report—a multi-layered cognitive model separating perception, inference, and abstraction—provides a robust and scalable foundation. Key to its success will be the deliberate focus on moving beyond simple correlation to robust causal inference, a step that circumvents the critical flaws of simpler imitation learning approaches. By modeling the game's rules rather than merely mimicking player actions, the system will generate knowledge that is both deeper and more generalizable.

### Actionable Recommendations:

1. **Prioritize the Unified JSON Schema:** The development of the detailed JSON schema (as outlined in Section 4.1) should be the immediate first step, undertaken collaboratively by the teams responsible for the analysis pipeline (Phase 2) and the Unity integration (Phase 3). This schema is the central contract of the project; finalizing it early will enable parallel development and prevent costly integration failures.
2. **Adopt the "Model as a Game" Consistency Framework:** The principles of numerical and spatial consistency using external state modules (Section 5.1) should be integrated into the architecture from the outset. Retrofitting these consistency checks later will be significantly more difficult. This ensures the pipeline produces reliable and logically sound outputs.
3. **Invest Heavily in Data Annotation and Management:** The performance of the Perception and Inference layers is directly proportional to the quality and quantity of the training data. A dedicated effort must be made to build the tooling and processes for annotating object bounding boxes, UI elements, and, most critically, time-stamped gameplay events. This annotated data is the fuel for the entire system.
4. **Develop in Stages with Validation at Each Step:** The pipeline should be built and validated incrementally.
  - **Stage 1:** Focus on the Perception Layer for a single, well-defined genre (e.g., a 2D platformer). Validate the accuracy of object detection, tracking, and UI extraction against a ground-truth game.
  - **Stage 2:** Build the Inference and Abstraction layers for that same genre. Validate the inferred rules and generated JSON against the game's known mechanics.
  - **Stage 3:** Once the pipeline is proven for one genre, begin scaling out to other genres by training new specialized models.

By following this structured, principled, and phased approach, the vision of reverse-engineering game design at scale can be transformed from an ambitious concept into a powerful, production-ready reality. This system will not only accelerate content creation but will also provide a unique, data-driven lens through which to understand the very art of game design itself.

## Works cited

1. Game Design Extractor.pdf
2. Procedural Content Generation for video games, a friendly approach, accessed June 24, 2025,  
<https://www.levelup-gamedevhub.com/en/news/procedural-content-generation->

- [for-video-games-a-friendly-approach/](#)
3. Procedural generation - Wikipedia, accessed June 24, 2025, [https://en.wikipedia.org/wiki/Procedural\\_generation](https://en.wikipedia.org/wiki/Procedural_generation)
  4. Mastering AI-Powered Procedural Content Generation for Games - Whimsy Games, accessed June 24, 2025, <https://whimsygames.co/blog/mastering-ai-powered-procedural-content-generation-for-games/>
  5. A Survey on Game Playing Agents and Large Models: Methods, Applications, and Challenges - arXiv, accessed June 24, 2025, <https://arxiv.org/html/2403.10249v1>
  6. Model as a Game: On Numerical and Spatial Consistency for Generative Games - arXiv, accessed June 24, 2025, <https://arxiv.org/html/2503.21172v1>
  7. Computer Vision in Gaming | Ultralytics, accessed June 24, 2025, <https://www.ultralytics.com/blog/computer-vision-in-gaming-leveling-up-player-experience>
  8. How to implement object detection for live streaming of games with fixed layout? - Reddit, accessed June 24, 2025, [https://www.reddit.com/r/computervision/comments/10amwh4/how\\_to\\_implement\\_object\\_detection\\_for\\_live/](https://www.reddit.com/r/computervision/comments/10amwh4/how_to_implement_object_detection_for_live/)
  9. Comprehensive Guide to Object Tracking in Computer Vision - Ikomia, accessed June 24, 2025, <https://www.ikomia.ai/blog/object-tracking-computer-vision-guide>
  10. Computer Vision for Object Detection and Tracking - Nexgits, accessed June 24, 2025, <https://nexgits.com/computer-vision-for-object-detection-and-tracking/>
  11. Activity Recognition from Video and Optical Flow Data Using Deep Learning - MathWorks, accessed June 24, 2025, <https://www.mathworks.com/help/vision/ug/activity-recognition-from-video-and-optical-flow-using-deep-learning.html>
  12. Visual Analysis of Computer Game Output Video Stream for ..., accessed June 24, 2025, <https://pmc.ncbi.nlm.nih.gov/articles/PMC7302570/>
  13. Action Recognition In Videos | Papers With Code, accessed June 24, 2025, <https://paperswithcode.com/task/action-recognition-in-videos-2>
  14. Gameplay highlight detection using Multi-modal learning, accessed June 24, 2025, <https://www.cs.toronto.edu/~edithal/2024-02-06-game-event-detection/>
  15. (PDF) Causal Inference with Observational Data - ResearchGate, accessed June 24, 2025, [https://www.researchgate.net/publication/24096771\\_Causal\\_Inference\\_with\\_Observational\\_Data](https://www.researchgate.net/publication/24096771_Causal_Inference_with_Observational_Data)
  16. What exactly is causal inference? How do you use it in your job? : r/datascience - Reddit, accessed June 24, 2025, [https://www.reddit.com/r/datascience/comments/18xtii1/what\\_exactly\\_is\\_causal\\_inference\\_how\\_do\\_you\\_use/](https://www.reddit.com/r/datascience/comments/18xtii1/what_exactly_is_causal_inference_how_do_you_use/)
  17. Pls help me learn causal inference from observational data : r/AskStatistics - Reddit, accessed June 24, 2025, [https://www.reddit.com/r/AskStatistics/comments/19dlnin/pls\\_help\\_me\\_learn\\_causal\\_inference\\_from/](https://www.reddit.com/r/AskStatistics/comments/19dlnin/pls_help_me_learn_causal_inference_from/)

18. Autumn: Causal Discovery Through Program Synthesis - Basis, accessed June 24, 2025, <https://www.basis.ai/blog/autumn/>
19. View of Mechanic Maker: Accessible Game Development via Symbolic Learning Program Synthesis - AAAI Publications, accessed June 24, 2025, <https://ojs.aaai.org/index.php/AIIDE/article/view/31884/34051>
20. Seconds to infer an FSM from 100 traces. | Download Table - ResearchGate, accessed June 24, 2025, [https://www.researchgate.net/figure/Seconds-to-infer-an-FSM-from-100-traces\\_tbl2\\_326335948](https://www.researchgate.net/figure/Seconds-to-infer-an-FSM-from-100-traces_tbl2_326335948)
21. Supervised and Unsupervised learning - GeeksforGeeks, accessed June 24, 2025, <https://www.geeksforgeeks.org/machine-learning/supervised-unsupervised-learning/>
22. Monte Carlo tree search - Wikipedia, accessed June 24, 2025, [https://en.wikipedia.org/wiki/Monte\\_Carlo\\_tree\\_search](https://en.wikipedia.org/wiki/Monte_Carlo_tree_search)
23. Behavioural cloning in video games, accessed June 24, 2025, <https://erepo.uef.fi/server/api/core/bitstreams/cfed9a16-0412-4be9-b797-d0f14e699155/content>
24. Behavioral Cloning in Atari Games Using a Combined Variational Autoencoder and Predictor Model, accessed June 24, 2025, [https://www.alexgorodetsky.com/static/papers/cheng\\_tandon\\_gorsich\\_gorodetsky\\_veerapaneni\\_ieee\\_ccc\\_2021.pdf](https://www.alexgorodetsky.com/static/papers/cheng_tandon_gorsich_gorodetsky_veerapaneni_ieee_ccc_2021.pdf)
25. Free AI Game Asset Generator | CGDream, accessed June 24, 2025, <https://cgdream.ai/features/ai-game-asset>
26. Meshy AI - The #1 AI 3D Model Generator for Creators, accessed June 24, 2025, <https://www.meshy.ai/>
27. 3D AI Studio - Generate 3D Models from Image or Text in Seconds, accessed June 24, 2025, <https://www.3daistudio.com/>
28. The Best Text to Game Asset AI Generator (for Free) - OpenArt, accessed June 24, 2025, <https://openart.ai/generator/game-asset>
29. Natural Language Generation (NLG) - Deepgram, accessed June 24, 2025, <https://deepgram.com/ai-glossary/natural-language-generation>
30. Natural Language Generation: Automate Text Creation with NLG - Lyzr AI, accessed June 24, 2025, <https://www.lyzr.ai/glossaries/natural-language-generation/>
31. Natural Language Generation Explained & 2 How To Tutorials In Python - Spot Intelligence, accessed June 24, 2025, <https://spotintelligence.com/2023/09/26/natural-language-generation/>
32. Prompt Engineering for AI Guide | Google Cloud, accessed June 24, 2025, <https://cloud.google.com/discover/what-is-prompt-engineering>
33. Prompt Engineering for Game Development - Analytics Vidhya, accessed June 24, 2025, <https://www.analyticsvidhya.com/blog/2024/06/prompt-engineering-for-game-development/>
34. Prompt engineering - Wikipedia, accessed June 24, 2025,

[https://en.wikipedia.org/wiki/Prompt\\_engineering](https://en.wikipedia.org/wiki/Prompt_engineering)

35. Model as a Game: On Numerical and Spatial Consistency for Generative Games, accessed June 24, 2025,  
[https://www.researchgate.net/publication/390247172\\_Model\\_as\\_a\\_Game\\_On\\_Numerical\\_and\\_Spatial\\_Consistency\\_for\\_Generative\\_Games](https://www.researchgate.net/publication/390247172_Model_as_a_Game_On_Numerical_and_Spatial_Consistency_for_Generative_Games)
36. [Literature Review] Model as a Game: On Numerical and Spatial Consistency for Generative Games - Moonlight | AI Colleague for Research Papers, accessed June 24, 2025,  
<https://www.themoonlight.io/en/review/model-as-a-game-on-numerical-and-spatial-consistency-for-generative-games>
37. Quantifying Video Gaming Expertise, accessed June 24, 2025,  
<https://bpb-us-e1.wpmucdn.com/sites.psu.edu/dist/7/89111/files/2020/12/2020-Elliott-Hampton-Quantifying-Video-Gaming-Expertise-preprint-a.pdf>