
To Infinity and Beyond: SHOW-1 and Showrunner Agents in Multi-Agent Simulations

Philipp Maas
Fable Studio

Frank Carey
Fable Studio

Chris Wheeler
Fable Studio

Edward Saatchi
Fable Studio

Pete Billington
Fable Studio

Jessica Yaffa Shamash
Fable Studio



Abstract

1 In this work we present our approach to generating high-quality episodic content for
2 IP's (Intellectual Property) using large language models (LLMs), custom state-of-
3 the-art diffusion models and our multi-agent simulation for contextualization, story
4 progression and behavioral control. Powerful LLMs such as GPT-4 were trained on
5 a large corpus of TV show data which lets us believe that with the right guidance
6 users will be able to rewrite entire seasons. *"That Is What Entertainment Will Look
7 Like. Maybe people are still upset about the last season of Game of Thrones.
8 Imagine if you could ask your A.I. to make a new ending that goes a different way
9 and maybe even put yourself in there as a main character or something."*¹

¹Brockman <https://www.hollywoodreporter.com/business/digital/chatgpt-game-of-thrones-openai-greg-brockman-1235348099/amp/>

10 1 Creative limitations of existing generative AI Systems

11 Current generative AI systems such as Stable Diffusion (Image Generator) and ChatGPT (Large
12 Language Model) excel at short-term general tasks through prompt engineering. However, they do
13 not provide contextual guidance or intentionality to either a user or an automated generative story
14 system (showrunner²) as part of a long-term creative process which is often essential to producing
15 high-quality creative works, especially in the context of existing IP's.

16 1.1 Living with uncertainty



Figure 1: Example Still from South Park AI Episode

17 By using a multi-agent³ simulation as part of the process we can make use of data points such as
18 a character's history, their goals and emotions, simulation events and localities to generate scenes
19 and image assets more coherently and consistently aligned with the IP story world. The IP-based
20 simulation also provides a clear, well known context to the user which allows them to judge the
21 generated story more easily. Moreover, by allowing them to exert behavioral control over agents,
22 observe their actions and engage in interactive conversations, the user's expectations and intentions
23 are formed which we then funnel into a simple prompt to kick off the generation process.

24 Our simulation is sufficiently complex and non-deterministic to favor a positive disconfirmation.
25 Amplification effects can help mitigate what we consider an undesired "slot machine" effect which
26 we'll briefly touch on later. We are used to watching episodes passively and the timespan between
27 input and "end of scene/episode" discourages immediate judgment by the user and as a result reduces
28 their desire to "retry". This disproportionality of the user's minimal input prompt and the resulting
29 high-quality long-form output in the form of a full episode is a key factor for positive disconfirmation.

30 While using and prompting a large language model as part of the process can introduce "several
31 challenges".⁴ Some of them, like hallucinations, which introduce uncertainty or in more creative
32 terms "unexpectedness", can be regarded as creative side-effects to influence the expected story
33 outcome in positive ways. As long as the randomness introduced by hallucination does not lead to
34 implausible plot or agent behavior and the system can recover, they act as happy-accidents⁵, a term
35 often used during the creative process, further enhancing the user experience.

²<https://fablesimulation.com/blog/friends-ai-sitcom-simulation>

³Sung Park <https://arxiv.org/abs/2304.03442>

⁴Li <https://arxiv.org/abs/2303.17760>

⁵Maas <https://noproscaenium.com/from-a-i-character-to-sundance-filmmaker-with-gpt-3-d4ab80c31b4e>

36 **1.2 The Issue of ‘The Slot Machine Effect’ in current Generative AI tools**

37 The Slot Machine Effect refers to a scenario where the *generation of AI-produced content feels more*
38 *like a random game of chance rather than a deliberate creative process*⁶. This is due to the often
39 unpredictable and instantaneous nature of the generation process.

40 Current off-the-shelf generative AI systems do not support or encourage multiple creative evaluation
41 steps in context of a long-term creative goal. Their interfaces generally feature various settings, such
42 as sliders and input fields which increase the level control and variability. The final output however,
43 is generated almost instantaneously by the press of a button. This instantaneous generation process
44 results in immediate gratification providing a dopamine rush to the user. This reward mechanism
45 would be generally helpful to sustain a multi-step creative process over long periods of time but
46 current interfaces, the frequency of the reward and a lack of progression (stuck in an infinite loop)
47 can lead to negative effects such as frustration, the intention-action gap⁷ or a *loss of control over the*
48 *creative process. The gap results from behavioral bias favoring immediate gratification*, which can
49 be detrimental to long-term creative goals.

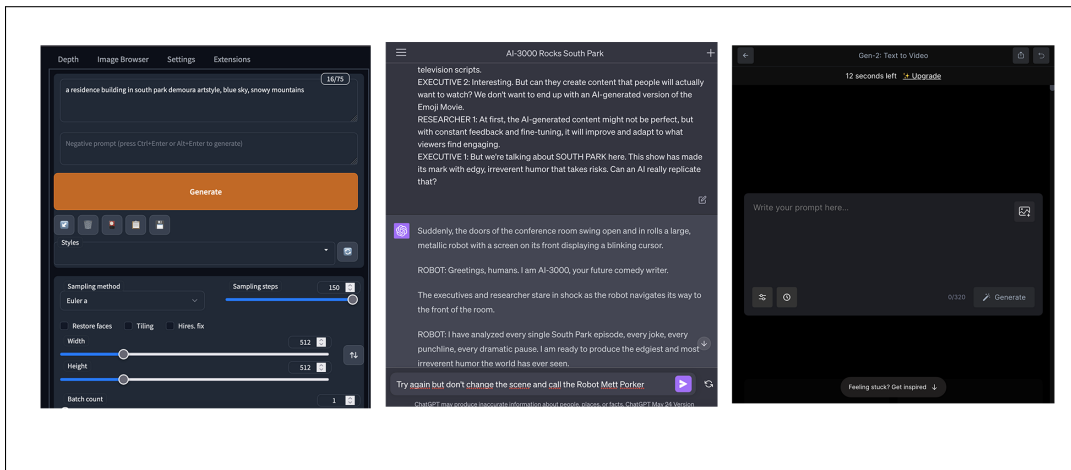


Figure 2: User Interface Comparison - Left to right: *Stable Diffusion Gradio*, *ChatGPT*, *Runway Gen-2*

50 While we do not directly solve these issues through interfaces, the contextualization of the process
51 in a simulation and the above mentioned disproportionality and timespan between input and output
52 help mitigate them. In addition we see opportunities in the simulation for in-character discriminators
53 that participate in the creative evaluation process, such as an agent reflecting on the role they were
54 assigned to or a scene they should perform in.

55 The multi-step "trial and error" process of the generative story system is not presented to the user,
56 therefore it doesn't allow for intervention or judgment, avoiding the negative effects of immediate
57 gratification through a user's "accept or reject" decisions. It does not matter to the user experience
58 how often the AI system has to retry different prompt chains⁸ as long as the generation process is
59 not negatively perceived as idle time but integrated seamlessly with the simulation gameplay. The
60 user only acts as the discriminator in the very end of the process after having watched the generated
61 scene or episode. This is also an opportunity to utilize the concept of Reinforcement Learning
62 through Human Feedback (RLHF) for improving the multi-step creative process and as a result the
63 automatically generated episode.

64 **1.3 Large Language Models**

65 LLMs represent the forefront of natural language processing and machine learning research, demon-
66 strating exceptional capabilities in understanding and generating human-like text. They are typically

⁶<https://artificial.tech/slot-machine-effect-of-ai/>

⁷<https://thedecisionlab.com/reference-guide/psychology/intention-action-gap>

⁸Yang <https://arxiv.org/abs/2306.02224>

67 built on Transformer-based architectures, a class of models that rely on self-attention⁹ mechanisms.
 68 Transformers allow for efficient use of computational resources, enabling the training of significantly
 69 larger language models. GPT-4, for instance, comprises billions of parameters that are trained on
 70 extensive datasets, effectively encoding a substantial quantity of worldly knowledge in their weights.

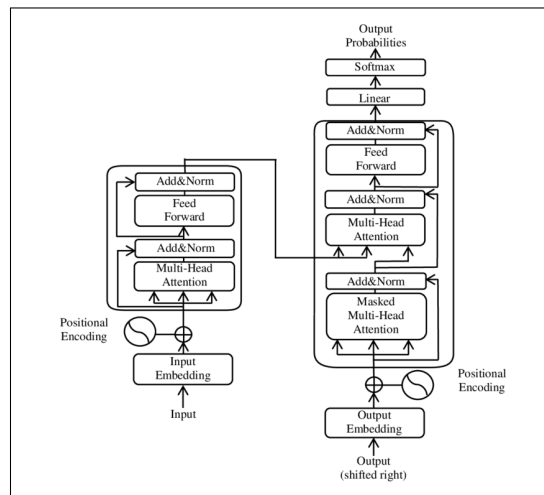


Figure 3: Diagram of the Transformer Architecture¹⁰

71 Central to the functioning of these LLMs is the concept of vector embeddings. These are mathematical
 72 representations of words or phrases in a high-dimensional space. These embeddings capture the
 73 semantic relationships between words, such that words with similar meanings are located close to
 74 each other in the embedding space. In the case of LLMs, each word in the model's vocabulary is
 75 initially represented as a dense vector, also known as an embedding. These vectors are adjusted
 76 during the training process, and their final values, or "embeddings", represent the learned relationships
 77 between words. During training, the model learns to predict the next word in a sentence by adjusting
 78 the embeddings and other parameters to minimize the difference between the predicted and actual
 79 words. The embeddings thus reflect the model's understanding of words and their context. Moreover,
 80 because Transformers can attend to any word in a sentence regardless of its position, the model
 81 can form a more comprehensive understanding of the meaning of a sentence. This is a significant
 82 advancement over older models that could only consider words in a limited window. The combination
 83 of vector embeddings and Transformer-based architectures in LLMs facilitates a deep and nuanced
 84 understanding of language, which is why these models can generate such high-quality, human-like
 85 text.

86 As was mentioned previously, transformer-based language models excel at short-term general tasks.
 87 They are regarded as fast-thinkers. [Kahneman]¹². Fast thinking pertains to instinctive, automatic,
 88 and often heuristic-based decision-making, while slow thinking involves deliberate, analytical, and
 89 effortful processes. LLMs generate responses swiftly based on patterns learned from training data,
 90 without the capacity for introspection or understanding the underlying logic behind their outputs.
 91 However, this also implies that LLMs lack the ability to deliberate, reason deeply, or learn from
 92 singular experiences¹³ in the way that slow-thinking entities, such as humans, can. While these
 93 models have made remarkable strides in text generation tasks, their fast-thinking nature may limit
 94 their potential in tasks requiring deep comprehension or flexible reasoning. More recent approaches
 95 to imitate slow-thinking capabilities such as prompt-chaining (see Auto-GPT) showed promising
 96 results. Large language models seem powerful enough to act as their own discriminator in a multi-step

⁹Vaswani <https://arxiv.org/abs/1706.03762>

¹⁰Vaswani <https://arxiv.org/abs/1706.03762>

¹¹<https://techcommunity.microsoft.com/t5/azure-data-explorer-blog/azure-data-explorer-for-vector-similarity-search/ba-p/3819626>

¹²Bubeck <https://arxiv.org/abs/2303.12712>

¹³Bubeck <https://arxiv.org/abs/2303.12712>

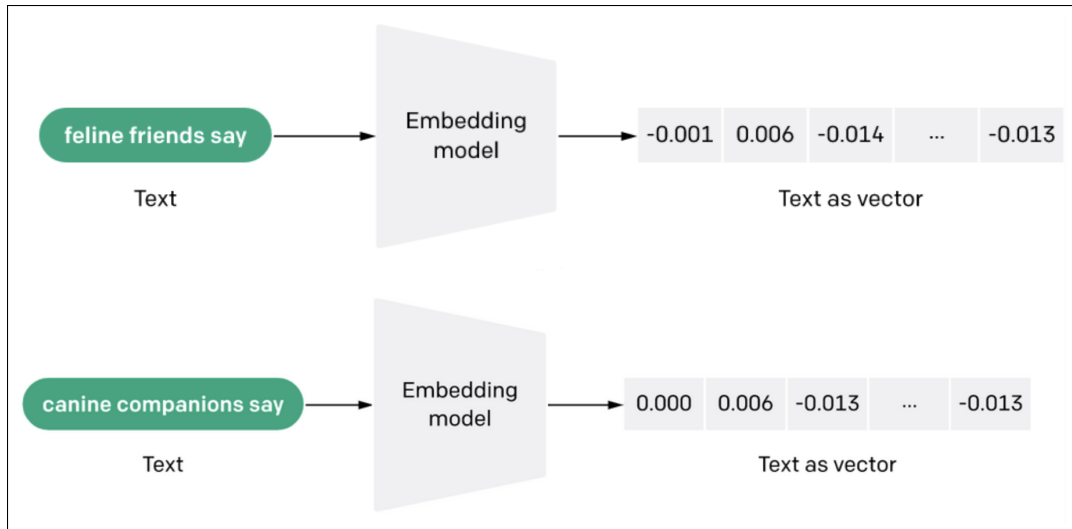


Figure 4: Example of Text Vector Embedding¹¹

97 process. *This can dramatically improve the ability to reason in different contexts, such as solving*
 98 *math problems.*¹⁴

99 We make heavy use of GPT-4 to influence the agents in the simulation as well as generating the
 100 scenes for the south park episode. Since transcriptions of most of the south park episodes are part of
 101 GPT-4’s training dataset, it already has a good understanding of the character’s personalities, talking
 102 style as well as overall humor of the show, eliminating the need for a custom fine-tuned model.

103 However, we do imitate slow-thinking as part of a multi-step creative process. For this we use
 104 different prompt chains to compare and evaluate the events of different scenes and how they progress
 105 the overall story towards a satisfactory, IP-aligned result. Our attempt to generate episodes through
 106 prompt-chaining is due to the fact that story generation is a highly discontinuous task.¹⁵ *These are*
 107 *tasks where the content generation cannot be done in a gradual or continuous way, but instead*
 108 *requires a certain "Eureka" idea that accounts for a discontinuous leap in the progress towards*
 109 *the solution of the task. The content generation involves discovering or inventing a new way of*
 110 *looking at or framing the problem, that enables the generation of the rest of the content. Examples*
 111 *of discontinuous tasks are solving a math problem that requires a novel or creative application of*
 112 *a formula, writing a joke or a riddle, coming up with a scientific hypothesis or a philosophical*
 113 *argument, or creating a new genre or style of writing.*

114 1.4 Diffusion Models

115 Diffusion models operate on the principle of gradually adding or removing random noise from data
 116 over time to generate or reconstruct an output. The image starts as random noise and, over many
 117 steps, gradually transforms into a coherent picture, or vice versa.

118 In order to train our custom diffusion models, we collected a comprehensive dataset comprising
 119 approximately 1200 characters and 600 background images from the TV show South Park. This
 120 dataset serves as the raw material from which our models learned the style of the show.

121 To train these models, we employ Dream Booth.¹⁶ The result of this training phase is the creation of
 122 two specialized diffusion models.

123 The first model is dedicated to generating single characters set against a keyable background color.
 124 This facilitates the extraction of the generated character for subsequent processing and animation,
 125 allowing us to seamlessly integrate newly generated characters into a variety of scenes and settings. h

¹⁴Baker <https://openai.com/research/improving-mathematical-reasoning-with-process-supervision>

¹⁵Bubeck <https://arxiv.org/abs/2303.12712>

¹⁶Ruiz <https://arxiv.org/abs/2208.12242>

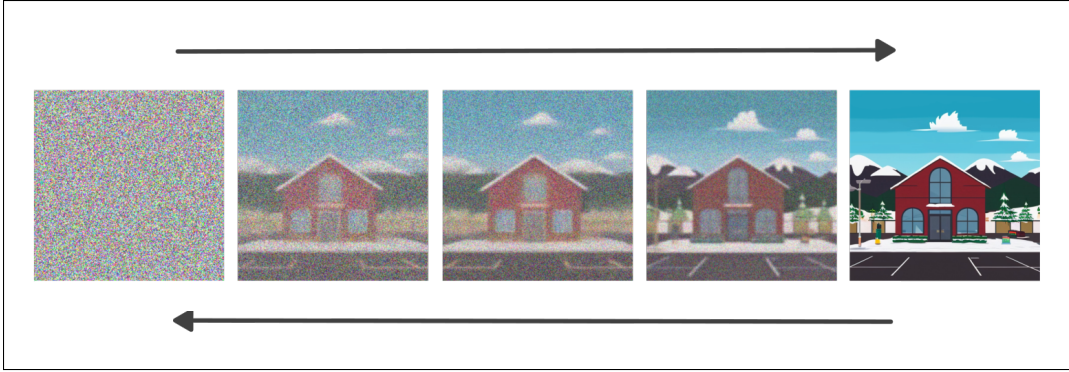


Figure 5: Stable Diffusion Model for South Park Backgrounds, prompt: *"a residence building in South park [demoura artstyle]"*

126

127 In addition, the character diffusion model allows the user to create a
 128 south park character based on their own looks via the image-to-image
 129 process of stable diffusion and then join the simulation as an equally
 130 participating agent. With the ability to clone their own voice, it's easy
 131 to imagine a fully realized autonomous character based on the user's
 132 characteristic looks, writing style and voice.



Figure 6: Example of generated South Park character

133 The second model is trained to generate clean backgrounds, with a
 134 particular focus on both exterior and interior environments. This model
 135 provides the 'stage' upon which our generated characters can interact,
 136 allowing for a wide range of potential scenes and scenarios to be created.



Figure 7: GPT-4 generated SVG image, prompt: *"Can you give me a svg drawing of a house on a street?"*

137

138 However, it's important to note that the images produced by these
 139 models are inherently limited in their resolution due to the pixel-based
 140 nature of the output. To circumvent this limitation, we post-process the
 141 generated images using an AI upscaling technique, specifically R-ESRGAN-4x+-Anime6B, which
 142 refines and enhances the image quality.

143 For future 2D interactive work, training custom transformer based models that are capable of
 144 generating vector-based output would have several advantages. Unlike pixel-based images, vector

145 graphics do not lose quality when resized or zoomed, thus offering the potential for infinite resolution.
 146 This will enable us to generate images that retain their quality and detail regardless of the scale at
 147 which they are viewed. Furthermore, vector based shapes are already separated into individual parts,
 148 solving pixel-based post-processing issues with transparency and segmentation which complicate the
 149 integration of generated assets into procedural world building and animation systems.

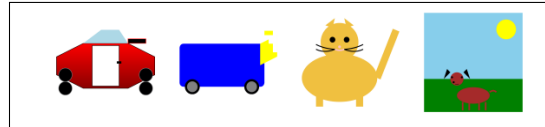


Figure 8: SVGs generated by GPT-4 for the classes automobile, truck, cat, dog.¹⁷

150 2 Episode Generation

151 We define an episode as a sequence of dialogue scenes in specific locations which add up to a total
 152 runtime of a regular 22 min south park episode.

153 In order to generate a full south park episode, we prompt the story system with a high level idea,
 154 usually in the form of a title, synopsis and major events we want to see happen over the course of 1
 155 week in simulation time (=roughly 3 hours of play time).

156 From this, the story system automatically extrapolates up to 14 scenes by making use of simulation
 157 data as part of a prompt chain. The showrunner system takes care of casting the characters for each
 158 scene and how the story should progress through a plot pattern. Each scene is associated with a plot
 159 letter (e.g. A, B, C) which is then used by the showrunner to alternate between different character
 groups and follow their individual storylines over the course of an episode to keep the user engaged.

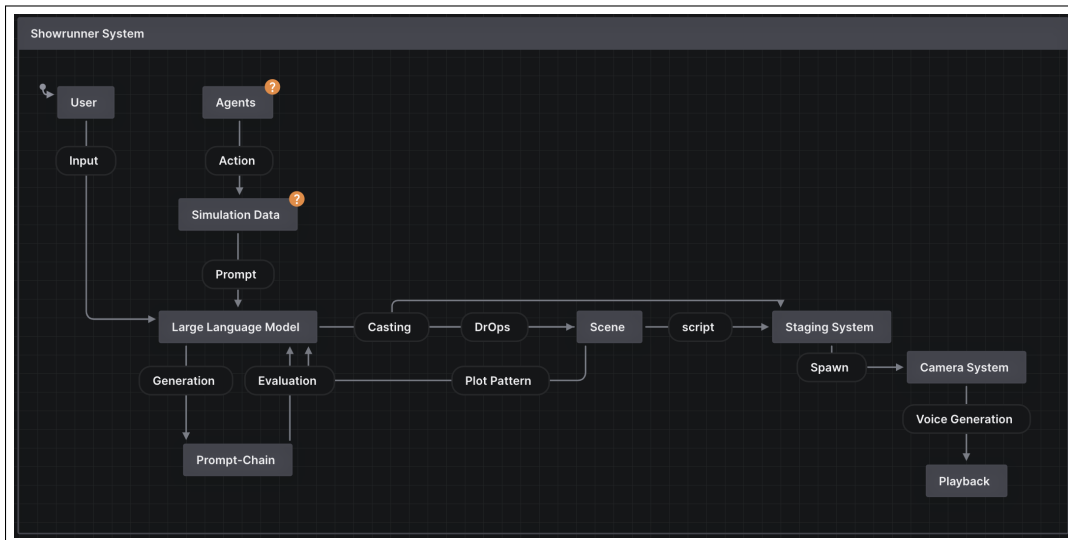


Figure 9: Diagram of Showrunner systems and prompt graph

160

161 In the end, each scene simply defines the location, cast and dialogue for each cast member. The scene
 162 is played back according to the plot pattern (e.g. ABABC) after the staging system and AI camera
 163 system went through initial setup. The voice of each character has been cloned in advance and voice
 164 clips are generated on the fly for every new dialogue line.

¹⁷<https://arxiv.org/abs/2303.12712>

165 **2.1 Reducing Latency**

166 In our experiments, generating a single scene can take a significant amount of time of up to one
167 minute. Below is a response time comparison between GPT-3.5-turbo and GPT-4. Speed will increase
168 in the short-term as models and service infrastructure get improved and other factors like artificial
169 throttling due to high user demand will get removed.

170 Since we generate the episodes during gameplay, we have ways to hide most of the generation time in
171 moments when the user is still interacting with the simulation or other user interfaces. Another way
172 to reduce the time needed to generate a scene or episode is to use faster models such as GPT-3.5-turbo
for specific prompts in the chain where the highest quality and accuracy is not so important.

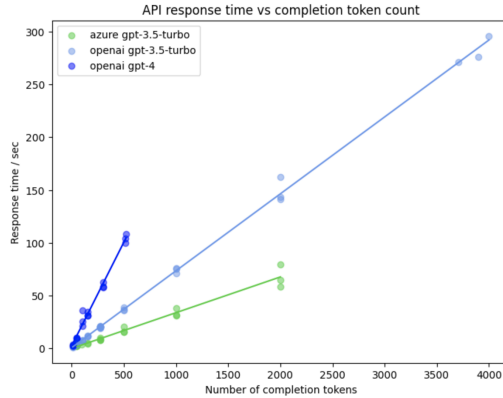


Figure 10: Speed comparison of GPT-3.5 vs. GPT-4¹⁸

173

174 During scene playback, we avoid any unwanted pauses between dialogue lines related to audio
175 generation by using a simple buffering system which generates at least one voice clip in advance.
176 See figure 11. This means while one character is delivering their voice clip, we already make the
177 web request for the next voice clip, wait for it to generate, download the file and then wait for the
178 current speaker to finish his dialogue before playback (delay). In this way the next dialogue line's
179 voice clip is always delivered without any delay. Text generation and voice cloning services become
180 increasingly fast and allow for highly adaptive and near-real time voice conversations.

181 **2.2 Simulate creative thinking**

182 As stated earlier, the data produced by the simulation acts as creative fuel to both, the user who is
183 writing the initial prompt and the generative story system which is interacting with the LLM via
184 prompt-chaining. Prompt-chaining¹⁹ is a technique, which involves supplying the language model

¹⁸Pungas <https://www.taivo.ai/gpt-3-5-and-gpt-4-response-times/>

¹⁹Wu <https://arxiv.org/abs/2203.06566>

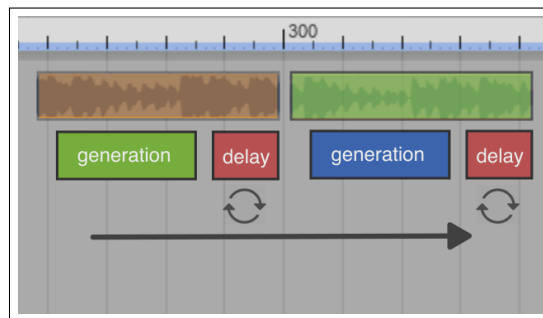


Figure 11: Diagram of zero-delay voice clip generation

185 with a sequence of related prompts to simulate a continuous thought process. Sometimes it can take
 186 on different roles in each step to act as the discriminator against the previous prompt and generated
 187 result.

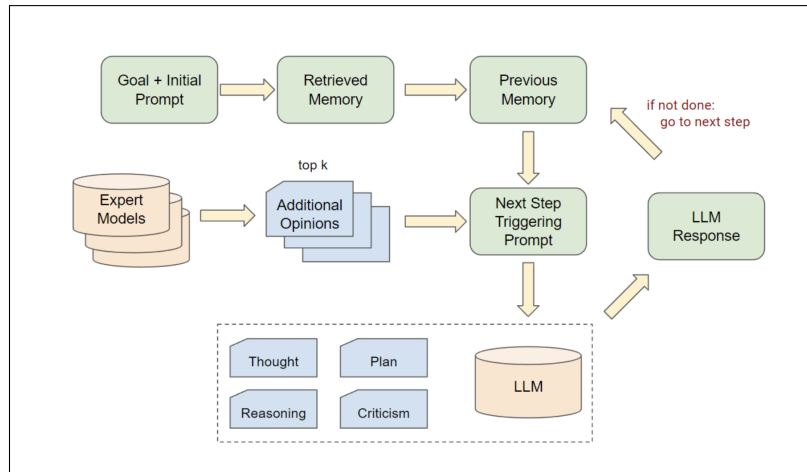


Figure 12: Example of a prompt chain from Auto-GPT²⁰

188 In our case we try to mimic that of a discontinuous creative thought process. For example, the
 189 creation of 14 distinct South Park scenes could be managed by initially providing a broad prompt to
 190 outline the general narrative, followed by specific prompts detailing and evaluating each scene’s cast,
 191 location, and key plot points. This mimics the process of human brainstorming, where ideas are built
 192 upon and refined in multiple often discontinuous steps. By leveraging the generative capabilities of
 193 LLMs in conjunction with the iterative refinement offered by prompt-chaining, we can effectively
 194 construct a dynamic, detailed, and engaging narrative.

195 In addition, we explore new concepts like plot patterns and dramatic operators (DrOps) to enhance
 196 the episode structure overall but also the connective tissue between each scene. Stylistic devices
 197 like reversals, foreshadowing, cliffhangers are difficult to evaluate as part of a prompt chain. A
 198 user without a writing background would have equal difficulty in judging these stylistic devices for
 199 their effectiveness and proper placement. We propose a procedural approach, injecting these show
 200 specific patterns and stylistic devices into the prompt chain programmatically as plot patterns and
 201 DrOps which can operate on the level of act structures, scene structures and individual dialogue
 202 lines. We are investigating future opportunities to extract what we call a dramatic fingerprint which
 203 is specific to each IP and format and train our custom SHOW-1 model with these data points. This
 204 dataset combined with overall human feedback could further align tone, style and entertainment value
 205 between the user and the specified IP while offering a highly adaptive and interactive story system as
 206 part of the on-going simulation.

207 2.3 Blank Page Problem

208 As mentioned above, one of the advantages of the simulation is that it avoids the blank page problem
 209 for both a user and a large language model by providing creative fuel²². Even experienced writers
 210 can sometimes feel overwhelmed when asked to come up with a title or story idea without any prior
 211 incubation of related material. The same could be said for LLMs. The simulation provides context
 212 and data points before starting the creative prompt chain.

²⁰Yank, Auto-GPT for Online Decision Making: Benchmarks and Additional Opinions
<https://arxiv.org/pdf/2306.02224.pdf>

²¹Drhlik, <https://pdrhlik.github.io/southparktalk-why2018/>

²²<https://www.trytriggers.com/blog-posts/overcoming-the-barrier-of-the-blank-page>

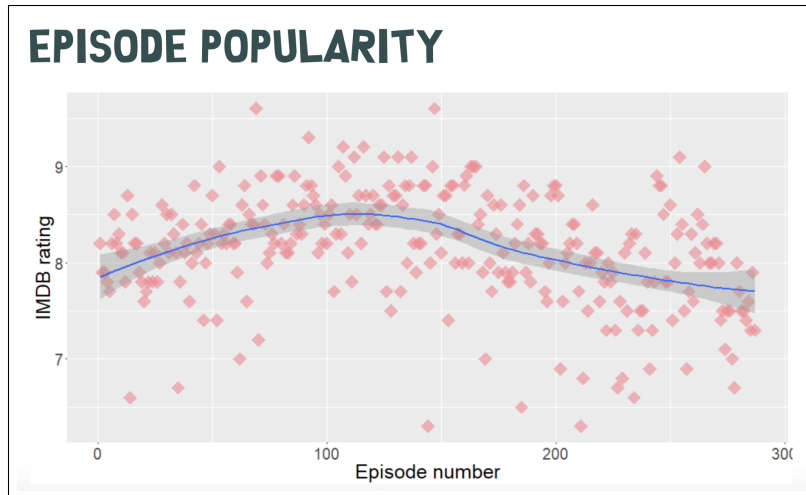


Figure 13: Diagram of South Park Episode ratings from IMDB²¹

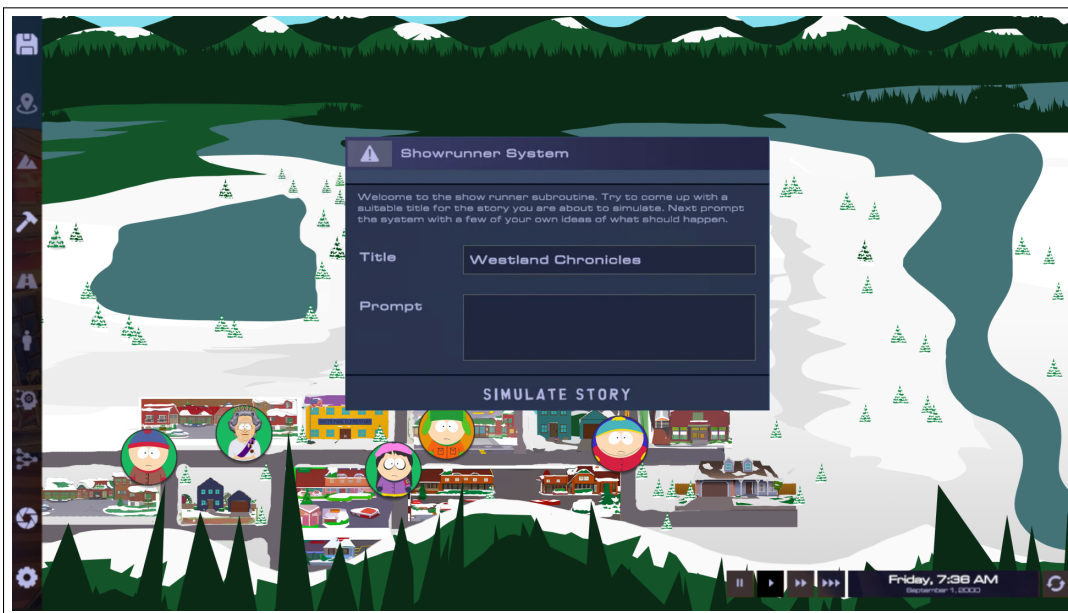


Figure 14: Example UI of Showrunner prompt input

213 **2.4 Who is driving the story**

214 The story generation process in our approach is a shared responsibility between the simulation,
 215 the user, and GPT-4. Each has strengths and weaknesses and a unique role to play depending on
 216 how much we want to involve them in the overall creative process. Their contributions can have
 217 different weights. While the simulation usually provides the foundational IP-based context, character
 218 histories, emotions, events, and localities that seed the initial creative process. The user introduces
 219 their intentionality, exerts behavioral control over the agents and provides the initial prompts that kick
 220 off the generative process. The user also serves as the final discriminator, evaluating the generated
 221 story content at the end of the process. GPT-4, on the other hand, serves as the main generative
 222 engine, creating and extrapolating the scenes and dialogue based on the prompts it receives from both
 223 the user and the simulation. It's a symbiotic process where the strengths of each participant contribute
 224 to a coherent, engaging story. Importantly, our multi-step approach in the form of a prompt-chain
 225 also provides checks and balances, mitigating the potential for unwanted randomness and allowing
 226 for more consistent alignment with the IP story world.

227 2.5 Intentionality

228 The formular (creative characteristics) and format (technical characteristics) of a show is often a
229 function of real-world limitations and production processes. They usually don't change, even over
230 the course of many seasons (South Park currently has 26 seasons and 325 episodes).²³

231 A single dramatic fingerprint of a show, which is used to train the proposed SHOW-1 model, can be
232 regarded as a highly variable template or "formula" for a procedural generator that produces South
233 Park-like episodes.

234 When a generative system is not limited in its ability to swiftly produce high amounts of content and
235 there is no limit for the user to consume such content immediately and potentially simultaneously,
236 the *10,000 Bowls of Oatmeal*²⁴ problem can become an issue. Everything starts to look and feel the
237 same or even worse, the user starts to recognize a pattern which in turn reduces their engagement as
238 they expect newly generated episodes to be like the ones before it, without any surprises.

239 This is quite different from a predictable plot which in combination with the above mentioned
240 "positive hallucinations" or happy accidents of a complex generative system can be a good thing.
241 Surprising the user by balancing and changing the phases of certainty vs. uncertainty helps increase
242 their overall engagement. If they would not expect or predict anything, they could also not get
243 pleasantly surprised.

244 With our work we aim for perceptual uniqueness. The OatMeal problem of procedural generators is
245 mitigated by making use of an on-going simulation (a hidden generator) and the long-form content of
246 22 min episodes which are only generated every 3h. This way the user generally does not consume a
247 high quantity of content simultaneously or in a very short amount of time. This artificial scarcity,
248 natural game play limits and simulation time help.

249 Another factor that keeps audiences engaged while watching a show and what makes episodes unique
250 is intentionality from the authors. A satirical moral premise, twisted social commentary, recent world
251 events or cameos by celebrities are major elements for South Park. Other show types, for example
252 sitcoms, usually progress mainly through changes in relationship (some of which are never fulfilled),
253 keeping the audience hooked despite following the same format and formula.

254 Intentionality from the user to generate a high-quality episode is another area of internal research.
255 Even users without a background in dramatic writing should be able to come up with stories, themes
256 or major dramatic questions they want to see played out within the simulation. To support this,
257 the showrunner system could guide the user by sharing its own creative thought process and make
258 encouraging suggestions or prompting the user by asking the right questions. A sort of reversed
259 prompt engineering where the user is answering questions.

260 One of the remaining unanswered questions in the context of intentionality is how much entertainment
261 value (or overall creative value) is directly attributed to the creative personas of living authors and
262 directors. Big names usually drive ticket sales but the creative credit the audience gives to the work
263 while consuming it seems different. Watching a Disney movie certainly carries with it a sense of
264 creative quality, regardless of famous voice actors, as a result of brand attachment and its history.

265 AI generated content is generally perceived as lower quality and the fact that it can get generated
266 in abundance further decreases its value. How much this perception would change if Disney were
267 to openly pride themselves on having produced a fully AI generated movie is hard to say. What if
268 Steven Spielberg, single handedly generated an AI movie? Our assumption is that the perceived value
269 of AI generated content would certainly increase.

270 A new interesting approach to replicate this could be the embodiment of creative AI models such
271 as SHOW-1 to allow them to build a persona outside their simulated world and build relationships
272 via social media²⁵ or real world events with their audience.²⁶ As long as an AI model is perceived
273 as a black box and does not share their creative process and reasoning in a human and accessible

²³<https://en.wikipedia.org/wiki/SouthPark>

²⁴Compton, Procedural Storytelling in Game Design

²⁵Virtual Beings <https://www.youtube.com/watch?v=FSq-mheA7Ds>, <https://www.youtube.com/watch?v=IROZSq-MQE>

²⁶Collaborating with AI at Sundance <https://www.fable-studio.com/behind-the-scenes/ai-collaboration>

274 way, as is the case for living writers and directors, it's unlikely to get credit with real creative values.
275 However, for now this is a more philosophical question in the context of AGI.

276 **3 Conclusion**

277 Our approach of using multi-agent simulation and large language models for generating high-quality
278 episodic content provides a novel and effective solution to many of the limitations of current AI
279 systems in creative storytelling. By integrating the strengths of the simulation, the user, and the
280 AI model, we provide a rich, interactive, and engaging storytelling experience that is consistently
281 aligned with the IP story world. Our method also mitigates issues such as the 'slot machine effect',
282 'the oatmeal problem' and 'blank page problem' that plague conventional generative AI systems.
283 As we continue to refine this approach, we are confident that we can further enhance the quality of
284 the generated content, the user experience, and the creative potential of generative AI systems in
285 storytelling.