

The Ultimate AI Experiment: When LLMs Play the Danganronpa Killing Game (Part 1)



ZACHARY HUANG

MAY 10, 2025



Ever wondered what happens when super-smart AI isn't just a tool, but a player in a crafty game of lies, strategy, and survival? What if AI could truly act out complex characters, scheme against each other, and fight for their lives in a twisted social deduction game? This two-part series shows you how we built exactly that: an AI-powered Danganronpa simulator where the characters act out the game themselves. You can even play along with the open-source code on GitHub.

Looks like an article worth saving!

Option

1. The LLM Could Unfulfilled Functionality

Remind me later

Hide Forever

Let's be real: Large Language Models (LLMs) are doing some mind-blowing stuff! – things that felt like pure sci-fi just a few years back! We've seen them crush super complex games like Go, write surprisingly human-like text and code, and even run simulations. AI agents can now explore and conquer virtual worlds like Minecraft (from NVIDIA's [Voyager](#)) and some can even beat entire RPGs like [Pokémon](#). The power is definitely there.

So, here's the big question: If LLMs are so smart, why aren't we seeing tons of awesome games where the AI characters (NPCs) actually feel *alive*? Why do many interactions in games still feel a bit... meh?

Thanks for reading Pocket Flow! Subscribe for free to receive new posts and support my work.

The Problem: NPCs Can Be Super Boring

We've all been there, right? You walk up to an NPC, and sure, they say their lines, they often lack that spark. They feel like slightly fancier versions of those stiff characters from older games – they're *there*, they *work*, but they rarely surprise you or make you feel truly connected. They just follow their scripts and don't really seem to have their own secret plans or motivations.

The result? Talking to them feels like a chore, and the game world, even with all its cool tech, can feel a bit lifeless. It feels like we're *this close* to LLMs creating truly awesome social experiences in games, but we're not quite there yet.

Looks like an article worth saving!

Option Q

Hover over the brain icon or use hotkeys to save with Memex.

Remind me later

Hide Forever



We think the secret sauce isn't just about more computer power. It's about giving AIs more personality and putting them in a game that forces them to show it. What if, instead of trying to build the perfect all-around NPC, we threw AIs into a crazy, high-stakes situation where they had to be more than just talking heads?

2. The Solution: Unleash AI into the Danganronpa Killing Game!

So, what did we do? We built **The Agentic Danganronpa Simulator**. Instead of making a new game from scratch, we used a game series famous for its wild characters and super-tense social gameplay: Danganronpa.

If you haven't heard of the twisted reality show where students (each an "Underdog" or "Star") get trapped by an escape room, a student has to murder a classmate and then get away with it in a "Class Trap".

Looks like an article worth saving!

Hover over the brain icon or use hotkeys to save with Memex.

[Remind me later](#) [Hide Forever](#)

In these trials, everyone debates and votes on who they think the killer (the "Blackened") is. Guess right? Only the killer is punished. Guess wrong? *Everyone* gets punished, and the killer goes free. Yikes!

Why Danganronpa? Because it's the perfect pressure cooker for AI!

This game is amazing for testing AI because it naturally has two things most AI experiences are missing:

- **Wild, Built-in Characters:** Danganronpa gives us a whole cast of "Ultimate" students with big personalities, unique ways of talking, and clashing goals. Imagine an AI trying to be Kokichi Oma, the "Ultimate Supreme Leader," who loves to lie and cause trouble. Or Miu Iruma, the "Ultimate Inventor," a super smart but foul-mouthed genius. These aren't just boring roles; they're like deep character studies waiting for an AI to jump in. The potential for crazy, character-driven moments is huge!
- **Gameplay That Creates Drama and Deception:** The whole "Killing Game" series is designed to make things intense. An AI can't just say it's a character; it has to act like them when things get real. Will the AI playing the nice guy Gonta Gokuhara freak out if he's falsely accused? How will the AI playing the smart Kyoko Kirigiri figure out the truth while keeping her own secrets? The constant danger, the need to lie, and the public accusations of the Class Trial force the AI to think strategically, team up, betray each other, and fight for their (virtual) lives.

Looks like an article worth saving!

Option Q

Hover over the brain icon or use hotkeys to save with Memex.

Remind me later

Hide Forever



Basically, we're not just asking LLMs to play Danganronpa. We're asking them to become characters. We want to see if their smarts and language skills can create the same kind of exciting, unpredictable, and very human (or inhuman!) drama that makes the original game awesome.

3. So, How Does This Twisted Game Actually Work?

Alright, picture this: 12 players are thrown into this high-stakes game of survival : deception. But here's the kicker – not everyone is who they seem!

Meet the Play

Looks like an article worth saving!

Option Q

Each player gets a se

Hover over the brain icon or use hotkeys to save with Memex.

- The Blackened (

Remind me later

Hide Forever

n?

eliminate other players without getting caught. The cool part is they know wh

their fellow Blackened are and secretly team up each night to choose one victim Diabolical!

- **The Truth-Seeker (Just 1!):** This player is like a secret detective. Each night, they can privately investigate one player to find out if they're a Blackened or an innocent Student. Crucial intel!
- **The Guardian (Also just 1!):** The protector! Each night, the Guardian can choose one player to shield from harm. But there's a catch: they can't pick the same player two nights in a row. Adds a bit of strategy, right?
- **The Students (The Innocent 7):** These are the regular folks just trying to survive. Their goal? To figure out who the Blackened are and vote them out during the Class Trials before it's too late.

The Daily Grind: Night, Morning, and... Trial!

The game unfolds in a few key phases:

1. Night Phase (When the Scheming Happens):

- First, the Blackened secretly chat amongst themselves and then vote on who to... *eliminate*. (Dun dun dun!)
- The Truth-Seeker picks one player to investigate, hoping to uncover a Blackened.
- The Guardian chooses someone to protect for the night.

2. Morning Phase (The Big Reveal):

- Monokuma pops up to announce what happened. Either a player has been eliminated (and their role is revealed – gasp!), or the Guardian successfully protected their target. Sometimes, the Blackened might even choose to abstain from voting.

3. Class Trial Phase

Looks like an article worth saving!

Option



Hover over the brain icon or use hotkeys to save with Memex.

- This only happens once per trial.

Remind me later

Hide Forever

01

a set order. This is where they present evidence, make accusations, or

desperately try to look innocent.

- **The Vote:** After the discussion, everyone **publicly** votes for who they think Blackened culprit is.
- **Execution (or not):** The player with the most votes gets kicked out, and the role is revealed. If there's a tie, nobody gets expelled, and the tension continues!

How Do You Win This Crazy Game?

There are two main teams, each with a different way to win:

- **Team Hope (Students, Truth-Seeker, Guardian):** They win if they successfully kick out all three Blackened. Justice prevails!
- **Team Despair (The Blackened):** They win if the number of Blackened becomes equal to or greater than the number of other living players. Basically, when despair takes over!

This setup is what makes the game a thrilling mix of deduction, bluffing, and social strategy. Seeing our AIs navigate this complex web of rules, roles, and relationships is what the Agentic Danganronpa Simulator is all about!

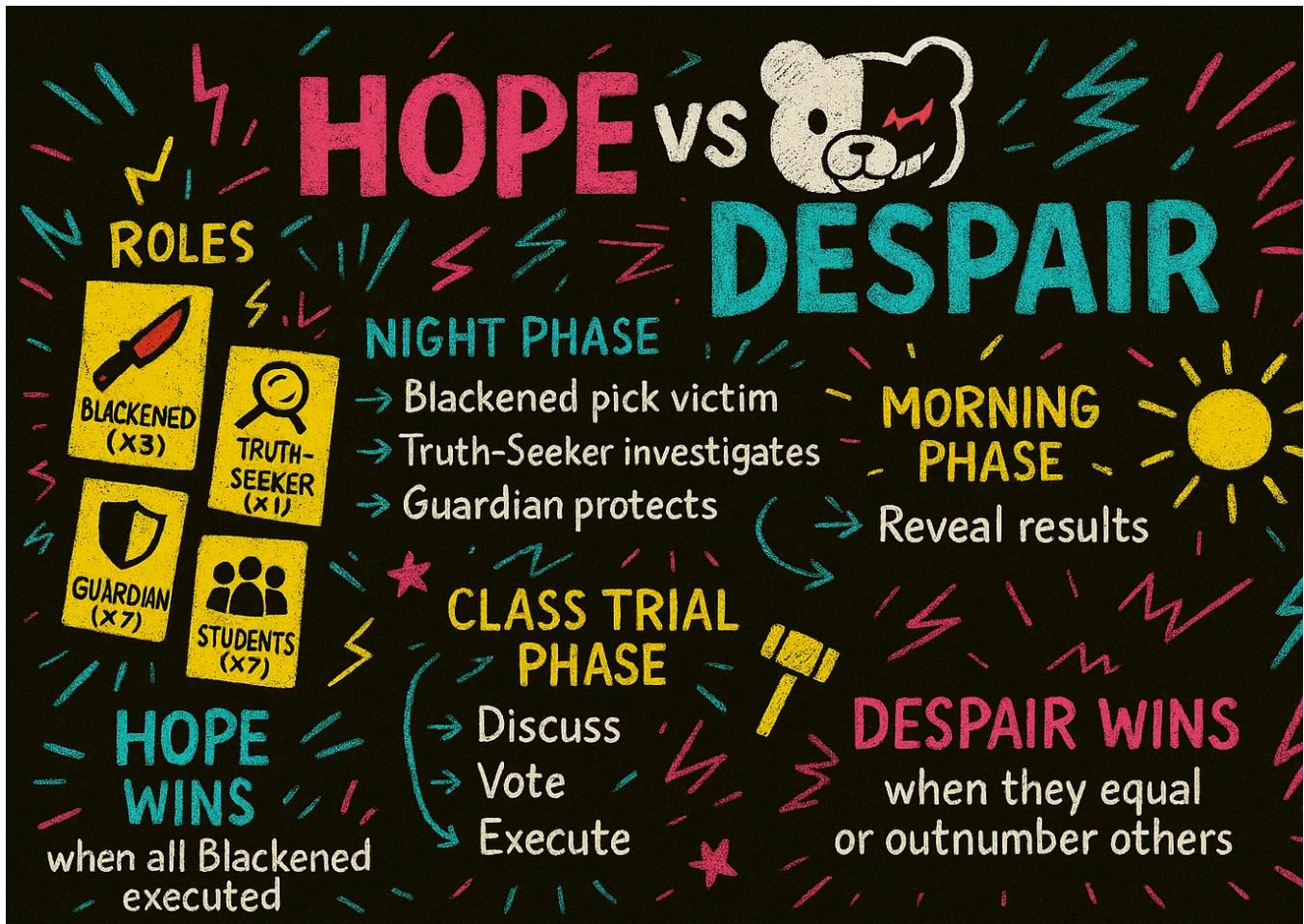
Looks like an article worth saving!

Option Q

Hover over the brain icon or use hotkeys to save with Memex.

Remind me later

Hide Forever



And guess what? You can experience this chaos firsthand:

- Play the Game NOW: <https://danganronpa.app/>
- Check Out the Code (It's Open Source!): <https://github.com/The-Pocket/PocketFlow-Tutorial-Danganronpa-Simulator>
- Watch Recorded Games (Seriously, It's Nuts!): <https://www.youtube.com/playlist?list=PLRYer4Da-4mJUSxS5oyn6GbXg4YWHe5Cr>

Now that you know how the game works, let's talk about what happened when we unleashed our AIs into it...

4. The AI Mi the Killer

Looks like an article worth saving!

Option Q



Hover over the brain icon or use hotkeys to save with Memex.

What really kicks the strategy and deception emerge. The AIs aren't just trying to figure out who

[Remind me later](#)

[Hide Forever](#)

so

they're deep in a cutthroat social game where persuading, manipulating, and even backstabbing are as crucial as actual clues.

It's fascinating to see these AI agents, powered by advanced language models, genuinely grapple with their assigned personalities and the game's pressures. They don't just recite lines; they *reason*, they *suspect*, they *panic*, and sometimes, they even *troll*, each approaching the challenges with its own information and biases.

It's a Social Battlefield, Not Just an Investigation:

- **The Persuasion Game (and the Lying Game!):** Knowing the truth is only half the battle. An AI Student who's fingered the Blackened still needs to *convince* other AIs. This involves crafting arguments, playing to others' personalities (or exploiting their weaknesses!), and navigating a minefield of deliberate misinformation spread by a cunning AI Blackened.
- **Lobbying and Alliance Building:** You'll see AIs trying to rally support, sometimes forming shaky voting blocs. A Blackened AI might subtly try to pin the blame on an innocent but unpopular character. A proactive Student AI might try to lead the charge, only to become a target for being too outspoken.
- **The Art of the Bluff:** When is it smart for an AI Truth-Seeker to reveal their research and findings? Doing so provides valuable information but also makes them a prime target. Some AIs might even *fake* being the Truth-Seeker – either to protect the real one or just to stir the pot (a classic Kokichi move!).
- **Calculated Betrayal:** What happens when an AI Blackened sees their teammates about to be exposed? Do they sink with them, or make the cold-blooded choice to "bus" their ally (throw them under to save themselves and look more trustworthy)? Our AIs have surprised us with their capacity for such ruthless (and strategically sound!) moves.

Freedom to Shape Strategy

Because the interactions

Looks like an article worth saving!

Option



Save

Hover over the brain icon or use hotkeys to save with Memex.

Remind me later

Hide Forever

strategic maneuvering. They aren't limited to predefined dialogue trees or act

menus. They can:

- **Filibuster or Misdirect:** An AI trying to buy time or protect an ally might lau into lengthy, irrelevant arguments.
- **Sow Seeds of Doubt:** Subtle insinuations or pointed questions can shift the group's focus, often based on flimsy (or entirely fabricated) "evidence."
- **Exploit Emotional States:** If an AI character is known to be gullible or prone panic, other AIs (especially a Blackened) might target them with manipulative tactics.

What kind of strategy would you devise if you were one of these AIs? Would you p safe, lay low, and observe? Or would you take bold risks, make daring accusations, try to control the flow of the Class Trial? The simulator allows these AI agents to explore a wide range of such tactics, often leading to emergent gameplay that feels remarkably organic.

Witnessing the "Inner Sanctum": The AI's Hidden Thoughts

One of the coolest parts is peeking into the AI's "thinking" process (which we log) a unique look at their reasoning and how they strategize:

- "Okay, [Character X] is accusing me, but their logic is flawed because [recap Y]. Best move: counter-accuse based on [detail Z], try to make them look flustered."
- "As a Blackened, my teammate [Character A] is drawing too much heat. If I don't voice them, we both go down. Sacrificing them makes me look like I'm helping the Students. Risky, but necessary."
- "I'm the Truth-Seeker; I know [Character B] is Blackened. But if I reveal myself now, other Blackened w ter without exposing i

Looks like an article worth saving!

Option Q

Hover over the brain icon or use hotkeys to save with Memex.

These internal mono
strategizing, reactin

Remind me later

Hide Forever

're
ali-

This depth turns the game from a simple simulation into a compelling psychological drama.



Ready to see some of this in action? Don't forget to check out the [Recorded Game Play](#) to witness these AI machinations firsthand!

5. The Code Behind the Chaos: Coming Up in Part 2!

We've seen how our Agentic Danganronpa Simulator can turn code and AI prompts into a surprisingly gripping show of personality, strategy, and betrayal. The AIs plot, accuse, defend, and cajole characters.

Looks like an article worth saving!

Option Q

on

Hover over the brain icon or use hotkeys to save with Memex.

But how does this all

Remind me later

Hide Forever

- How do we keep a wild Danganronpa game flowing smoothly, from secret nig actions to chaotic class trials, making sure every AI gets its turn and the right info?
- What does the tech setup look like to manage all these different AI agents?
- How do we make sure an AI playing a super-secretive character like Kokichi doesn't accidentally see the Truth-Seeker's private notes? (Spoiler: It involves some clever tricks with information control!)
- What kind of digital notebooks do we need to keep track of everyone's roles, moves, their "secret thoughts," and the constantly changing game state?
- And what about those little "game design" tweaks, like giving the main character a bit of "plot armor" so you (if you play as Shuichi) don't just get knocked out the first round every time by super-strategic AIs?

These are the juicy questions we'll dig into in [Part 2: Architecting Despair - Inside the Danganronpa Simulator](#). We'll pull back the curtain and get into the techy details exploring how the system is built, the game's digital flowchart, the key ways we handle secret information, how we log every whisper and accusation, and the tricks that keep the despair running smoothly.

If you're curious about building your own complex, AI-driven games, or just want to know the engineering secrets behind making AIs convincingly scheme and strategize, you won't want to miss it!

For now, why not [try to survive the Killing Game yourself](#), or [explore the open-source code on GitHub](#) and see if you can spot the beginnings of the chaos we've talked about? The despair is just getting started!

Looks like an article worth saving!

Option Q

Hover over the brain icon or use hotkeys to save with Memex.

Remind me later

Hide Forever



4 Likes

← Previous

Next

Discussion about this post

[Comments](#) [Restacks](#)

Write a comment...

© 2025 Zachary Huang · [Privacy](#) · [Terms](#) · [Collection notice](#)
[Substack](#) is the home for great culture

Looks like an article worth saving!

Option Q

Hover over the brain icon or use hotkeys to save with Memex.

Remind me later

Hide Forever