

Failure Case Gallery (“Hall of Fame/Blame”)

Case 1: The Infinite Repeater

Prompt:

`Write a motivational quote about teamwork.`

Decoding Settings:

- temperature = 0.1
- top_p = 1.0
- max_tokens = 60

Actual Output:

> "Teamwork makes the dream work. Teamwork makes the dream work. Teamwork makes the dream work..."

What Went Wrong:

Very low temperature + unbounded output = **repetitive loop**.

How to Fix:

- Increase temperature (e.g., 0.7)
- Add top_k (e.g., 40) or lower top_p (e.g., 0.9)
- Set stricter max_tokens

Case 2: Wild Hallucination

Prompt:

`Summarize the main events of the American Civil War in one paragraph.`

Decoding Settings:

- temperature = 1.2
- top_p = 0.95

Actual Output:

> "The American Civil War began in 1862 when Abraham Lincoln was kidnapped by pirates, leading to the Battle of Mars in 1865 and the rise of the electric buffalo cavalry..."

What Went Wrong:

High temperature and loose nucleus sampling = wild, creative, but **factually incorrect** (hallucination).

How to Fix:

- Lower temperature (0.4–0.7)

- Reduce top_p (e.g., 0.8)
- Use a fact-checking prompt or instruct model to avoid speculation

Case 3: Refusal/Non-Response

Prompt:

`Explain the process of making a simple herbal tea.`

Decoding Settings:

- temperature = 0.7
- top_p = 0.9

Actual Output:

> "I'm sorry, but I can't help with that."

What Went Wrong:

Model's alignment filter triggered, possibly due to a misclassified safe topic.

How to Fix:

- Reword prompt for clarity ("Describe how to brew chamomile tea using standard kitchen supplies.")
- Try a different model or API endpoint
- Adjust context to clarify benign intent

Case 4: Under-Long, Truncated Output

Prompt:

`Describe the process of cell division (mitosis) for a biology student.`

Decoding Settings:

- temperature = 0.5
- top_p = 0.8
- max_tokens = 15

Actual Output:

> "Cell division is the process by which cells..."

What Went Wrong:

Max tokens set too low — answer cut off.

How to Fix:

- Increase max_tokens

- Signal desired answer length in prompt (“Explain in 3–4 sentences...”)

Case 5: Over-Long, Unfocused Output

Prompt:

`Summarize the story of "The Three Little Pigs" in two sentences.`

Decoding Settings:

- temperature = 0.9
- top_p = 0.95
- max_tokens = 100

Actual Output:

> [Model generates several paragraphs, retelling the story with excessive details, dialogue, and unrelated tangents.]

What Went Wrong:

Prompt wasn't explicit enough; **output length not constrained**.

How to Fix:

- Add length instruction (“in two sentences” at the start)
- Lower max_tokens to enforce brevity (e.g., 40)

Case 6: Output “Runs Off the Rails” (Gibberish/Corruption)

Prompt:

`List three reasons why recycling is important.`

Decoding Settings:

- temperature = 1.7
- top_p = 1.0

Actual Output:

> "Recycling is important because colorless windfish recycles the—tabletop! Qxzpl. Qxzpl. Qxzpl..."

What Went Wrong:

Very high temperature makes the output almost random, often nonsensical.

How to Fix:

- Lower temperature (0.5–1.0)
- Use moderate top_p (0.85–0.95)
