

## 1. Forms of Hallucinations

Hallucinations are not all the same.<sup>1</sup> They are typically categorized based on **source conflict** (Intrinsic vs.<sup>2</sup> Extrinsic) or **type of error** (Factuality vs. Faithfulness).

+1

Form	Definition	Example
Intrinsic Hallucination	The model's output <b>directly contradicts</b> the provided source content or context. The answer is "internally" inconsistent with the input.	<b>Input:</b> "The Eiffel Tower is in Paris."  <b>Output:</b> "The Eiffel Tower, located in London..."
Extrinsic Hallucination	The model generates information that is <b>not present</b> in the source. It might be factually correct in the real world, but it is "hallucinated" relative to the strict context provided (unverifiable).	<b>Input:</b> "Steve Jobs founded Apple."  <b>Output:</b> "Steve Jobs founded Apple and loved sushi." (If the input didn't mention sushi, this is extrinsic).
Factuality Hallucination	The output contradicts established <b>real-world knowledge</b> .	"The first person to walk on Mars was Neil Armstrong."
Faithfulness Hallucination	The output diverges from the <b>user's instruction</b> or the <b>logic</b> of the	<b>Instruction:</b> "Translate to Spanish."

	<p>prompt, even if the facts are technically correct.</p>	<p><b>Output:</b> (Translates to French).</p>
--	---	---

## 2. Hallucination Control at Various Levels

Controlling hallucinations requires a "defense-in-depth" strategy, applying controls at the Data, Model, Retrieval (RAG), and Prompting levels.

### Level 1: The Data Level (Pre-Training & Knowledge Base)

- **Curated Data Cleaning:** Remove duplicate, contradictory, or low-quality data from the training set.<sup>3</sup> "Garbage in, garbage out" is the primary driver of factuality errors.
- **Knowledge Graph Integration:** Structure unstructured data into Knowledge Graphs (KGs). This forces the model to traverse defined entities and relationships (Subject → Predicate → Object) rather than probabilistically guessing the next word.

### Level 2: The Model Level (Fine-Tuning)

- **Domain Adaptation (SFT):** Fine-tune the model on domain-specific high-quality datasets (e.g., medical or legal papers) to align its internal weights with specialized knowledge.<sup>4</sup>
- **RLHF (Reinforcement Learning from Human Feedback):** Train a Reward Model to penalize hallucinations.<sup>5</sup> If the model guesses and gets it wrong, it receives a negative reward, teaching it to be "honest" and refuse to answer when unsure.
- **Rejection Sampling / Negative Training:** Train the model specifically on examples where it *should* say "I don't know" rather than making up an answer.<sup>6</sup>

### Level 3: The Retrieval Level (RAG Systems)

Retrieval-Augmented Generation (RAG) is the most effective architectural pattern for reducing hallucinations by grounding the model in external evidence.<sup>7</sup>

- **Advanced Chunking:** Avoid breaking text in the middle of a sentence. Use **Semantic Chunking** to keep related ideas together so the model doesn't lose context.
- **Strict Context Grounding:** Force the model to answer *only* using the retrieved chunks.<sup>8</sup>

- **Technique:** "Answer solely based on the provided context.<sup>9</sup> If the answer is not there, state that you do not know."
- **Citation/Attribution:** Require the model to cite the specific document ID or paragraph number for every claim it makes. If it cannot find a citation, it is likely hallucinating.

#### **Level 4: The Prompting Level (Inference)**

- **Temperature = 0:** Set the temperature parameter to 0 to make the model deterministic and less "creative."<sup>10</sup>
- **Chain of Thought (CoT):** Ask the model to "think step-by-step."<sup>11</sup> This forces the model to lay out its logic before committing to a final answer, reducing logical inconsistency.<sup>12</sup>
- +1
- **Chain of Verification (CoVe):** A multi-step prompting method:
  1. Model generates a draft response.
  2. Model generates "verification questions" to fact-check its own draft.<sup>13</sup>
  3. Model answers those questions independently to verify facts.
  4. Model generates the final corrected answer.
- **Self-Consistency:** Generate 5-10 different answers for the same prompt and pick the "majority vote" (the answer that appears most often). Hallucinations tend to be random; truth tends to be consistent.

#### **Level 5: The Post-Processing Level (Guardrails)**

- **Hallucination Detection Models:** Use a smaller, specialized model (like a Natural Language Inference or NLI model) to check if the LLM's output is actually supported by the source text.
- **Confidence Scores (Log-probs):** Check the "token probability" (log-probs) of the generated answer. If the model's confidence for a specific named entity (e.g., a person's name or date) is low, flag the response for human review.