

Architectural Comparison: Standard AI/RAG vs COMPAiSS

Standard Generation-First AI / RAG Architecture			COMPAiSS Execution-Gated Architecture		
STAGE	WHAT HAPPENS AT THIS STAGE	WHAT THIS MEANS FOR ACCURACY & RISK	STAGE	WHAT HAPPENS AT THIS STAGE	WHAT THIS MEANS FOR ACCURACY & RISK
Model Inference (Always Open)	The AI model is active and ready to generate responses by default. Inference begins as soon as a question is received.	<i>The system assumes it will provide an answer before knowing if one is appropriate.</i>	Model Inference (Always Open)	The AI model is not active by default. Inference is gated and will only occur if authorization succeeds.	<i>The system does not assume it will answer before determining if it should.</i>
User Question Intake	The user's question is received in any language and passed into the system for processing.	<i>No institutional scope or authorization check occurs at this point.</i>	User Question Intake	The user's question is received in any language and immediately evaluated for institutional scope and authorization.	<i>Questions outside institutional authority are blocked before any AI reasoning begins.</i>
Language Translation / Normalization	If the question is not in English, it may be translated as part of the retrieval and reasoning process.	<i>Translation can influence which documents are retrieved and how meaning is interpreted.</i>	Language Translation / Normalization	If the question is not in English, it is translated solely to normalize meaning. This is a linguistic step only, not a retrieval or reasoning step.	<i>Translation does not influence which sources are considered or how answers are determined.</i>
Embedding Generation	The question is converted into a mathematical representation to enable similarity search against stored documents.	<i>The system searches broadly across all available data without institutional filtering.</i>	Embedding Generation	The normalized question is converted into a mathematical representation, but only within a pre-bounded institutional truth space.	<i>Search is restricted to greenlisted, institution-approved sources before any AI reasoning occurs.</i>
Vector / Document Retrieval (RAG)	The system retrieves documents or passages that seem relevant based on similarity matching. This may include external, outdated, or non-authoritative sources.	<i>Retrieved content may not be institutionally verified or authorized.</i>	Vector / Document Retrieval (RAG)	The system retrieves only from authorized institutional sources. If no authoritative material exists, the process stops here.	<i>Only verified, institution-approved content can inform an answer.</i>
Prompt Assembly	The AI combines the user's question with retrieved documents to construct a prompt for generating an answer.	<i>The AI may fill gaps, infer context, or draw from training data beyond what was retrieved.</i>	Prompt Assembly	If authorization succeeds, the AI is instantiated and combines the question with authorized materials. The AI cannot infer beyond what is explicitly provided.	<i>The AI is confined to institutional sources and cannot fill gaps from general training data.</i>
Post-Generation Controls	After the answer is generated, it may pass through filters, safety checks, moderation layers, or human review to catch errors.	<i>Controls reduce but do not eliminate hallucinations or unsupported claims.</i>	Post-Generation Controls	Because the answer is based solely on authorized sources, post-generation filtering is unnecessary for accuracy or authority.	<i>Hallucinations based on unauthorized sources are structurally prevented, not retroactively caught.</i>
Final Answer Delivered	The answer is translated back (if needed) and presented to the user, often with disclaimers about potential inaccuracy.	<i>Users receive answers that may include unauthorized links, unverified information, or mixed sources.</i>	Final Answer Delivered	The answer is translated back (if needed) and presented to the user, grounded exclusively in institutional materials. If no answer was authorized, the user is informed.	<i>Users receive only precise, authorized answers or a clear statement that no answer is available.</i>

In this model, AI is open from the start and answers are controlled through stages, but hallucinations are not eliminated.

In this model, AI reasoning is gated before it begins, preventing unsupported answers from being generated at all.