



# Policy-Enforced RAG for HIPAA-Compliant Healthcare AI

# Table of Contents

**EXECUTIVE SUMMARY.....3**

**INTRODUCTION.....3**

**BACKGROUND: THE NEED FOR POLICY-GROUNDED RAG .....4**

**WHAT IS POLICY GROUNDED RAG .....5**

**WHAT MAKES IT POLICY-GROUNDED?.....6**

**WHY POLICY-GROUNDED RAG MATTERS IN CLAIMS SYSTEMS ..... 10**

**POLICY GROUNDED RAG FRAMEWORK..... 11**

**TECHNICAL IMPLEMENTATION BLUEPRINT .....17**

**SAMPLE POLICY GROUNDED RAG WORKFLOW FOR AN HEALTHCARE CLAIM REVIEW .....19**

**BENEFITS OF POLICY GROUNDED RAG ..... 20**

**TRADITIONAL VS POLICY-GROUNDED RAG ..... 23**

**USE CASES BEYOND CLAIMS ..... 24**

**CONCLUSION ..... 24**

## Executive summary

Policy-Grounded Retrieval-Augmented Generation (RAG) is an enterprise-grade approach to Generative AI that ensures all model-generated outputs are strictly anchored to official organizational policies, rules, and guidelines. In regulated industries—particularly U.S. healthcare—GenAI systems must operate within precise compliance boundaries. Policy-Grounded RAG enforces these boundaries by requiring that every AI-generated answer be supported by verifiable policy evidence.

In this whitepaper, we use healthcare as the primary example to illustrate how Policy-Grounded Retrieval-Augmented Generation (RAG) ensures safe, compliant, and accurate AI-driven decision support. The healthcare claims environment provides the clearest demonstration of the risks and the value of grounding AI strictly in approved documentation.

## Introduction

Modern enterprises increasingly depend on AI to support complex, rule-driven decisions. In sectors such as healthcare, insurance or any other domain that are ruled by Regulations/Policies caused by AI hallucination can lead to:

- Regulatory violations
- Incorrect claim decisions
- Inaccurate denials or approvals
- Member dissatisfaction
- Legal and financial exposure

Policy-Grounded RAG solves this by binding AI outputs to the organization's approved documentation corpus. Instead of allowing a model to reason freely, we enforce strict reliance on retrieved, sanctioned policy text.

## Background: The Need for Policy-Grounded RAG

Enterprises across healthcare, insurance, financial services, and other regulated industries increasingly rely on AI to interpret complex documentation and support operational decision-making. However, traditional AI and even standard Retrieval-Augmented Generation (RAG) systems were never designed to operate within strict regulatory frameworks. They often pull information from mixed sources, apply generalized reasoning, or generate answers not fully supported by authoritative policy text.

In environments where decisions must strictly follow defined rules—such as CMS guidelines, benefit policies, prior authorization rules, insurance coverage criteria, and contractual obligations—these shortcomings create substantial risk. Even minor deviations from official policy can lead to:

- Non-compliant determinations that violate federal or state regulations
- Improper denials or approvals affecting claim accuracy and reimbursement
- Costly provider disputes, appeals, and rework
- Audit exposure under CMS, NCQA, URAC, or external regulators
- Operational inefficiencies from inconsistent policy interpretation across staff

Traditional RAG cannot reliably ensure that model outputs remain fully aligned with approved documents. It retrieves text but still allows the model to reason freely or introduce unsupported logic, which can be unacceptable in policy-driven workflows.

This gap created the need for Policy-Grounded RAG—a rigorously constrained GenAI framework that forces the model to rely exclusively on curated, up-to-date policy materials. By grounding every answer in verifiable policy evidence and enforcing strict guardrails, Policy-Grounded RAG provides a safe, auditable, and compliant foundation for AI-enabled decision support in highly regulated environments.

## What is Policy Grounded RAG

Policy-grounded RAG (Retrieval-Augmented Generation) is a design pattern where the LLM (Large Language Model) is only allowed/forced to answer questions using authoritative internal documents. This helps in the domains where Regulations & Policies is a must to be followed.

This approach combines two layers. **Retrieval Layer** and **Generation layer**

For Instance, in a US Healthcare system

- **Retrieval Layer:**

Before the LLM answers, it searches your internal knowledge base:

- Medical policies
- Insurance coverage guidelines
- benefit manuals
- internal SOPs
- CMS Guidelines
- LOAs (letters of agreement), etc.

- **Generation Layer:**

The LLM uses retrieved text to form a response with citations, such as:

"This claim is not eligible for coverage because Policy CP-102, Section 3.2 requires prior authorization for this service."

The model is disallowed from inventing rules or relying on general knowledge. It answers only with what's explicitly present in the documents. So instead of "making things up," the model is answering with documents in hand. This avoids LLM/SLM Hallucination

## What makes it policy-grounded?

Policy-grounded” means you add strict constraints so the LLM: Policy grounding is enforced through four pillars:

### 1. Curated, Controlled Corpus

Only approved, up-to-date policy documents are indexed. No external web data. No drafts. No outdated PDFs.

### 2. Metadata-Driven Retrieval

*Documents are tagged by:*

- Plan and product type
- Line of business
- Effective dates
- Region or state applicability
- Medical category (e.g., imaging, cardiology)

### 3. Guard railed Prompting

Guard railed prompting is the core mechanism that ensures the LLM does not operate freely or creatively. Instead, it enforces strict policy compliance, factual grounding, and controlled behavior.

To make Policy-Grounded RAG safe for regulated sectors such as U.S. healthcare claims, prompting must:

- Constrain the model’s reasoning to retrieved policy text only.
- Avoid hallucinations by forbidding the introduction of new rules.
- Force transparency via required citations.
- Require refusal when the retrieved information is insufficient.
- Maintain auditability by producing consistent, reference-linked outputs.

*The LLM is instructed:*

- ***Use only Retrieved Documents***

The LLM is given explicit system instructions that all reasoning must originate only from retrieved text. This prevents the model from drawing on general medical knowledge or prior training data.

- "You may not use external medical knowledge."
- "If the required policy text was not retrieved, you must decline to answer."

- ***Required Uncertainty Handling***

When policy information is incomplete or missing, the AI must respond with a safe fallback.

- "If the retrieved policy text does not clearly address the question, respond with: 'Insufficient information in the retrieved policy documents to make a determination.'"

This ensures the AI never fabricates rules or coverage criteria.

- ***Mandatory Policy Citations***

Every output must reference the exact policy ID and section used.

- Ensures auditability for CMS, URAC, NCQA.
- Supports traceability in appeals.
- Prevents ambiguous or unverifiable explanations.

***Examples:***

- "Based on Policy CP-102, Section 4.1..."
- "Policy MP-17, Section 2.3 states..."

- **Structured Output Schema**

Prompts may include a required format:

```
{  
  
  "determination": "Decision/Recommendation made based on output",  
  
  "rationale": "text grounded in retrieved policy",  
  
  "citations": ["Reference sections"]  
  
}
```

This ensures consistency across all reviewers and workflows.

A Few sample prompts for Healthcare to illustrate the Prompting of Policy-Grounded RAG

### **Sample Guard-railed Prompts (Healthcare Claims Sector)**

Ex 1. Sample Prompt for a Medical coverage review workflow

*"You are a Policy-Grounded Medical Coverage Assistant.  
Use ONLY the retrieved policy text shown below. Do NOT rely on external knowledge.  
If the retrieved text is insufficient, respond with: "Insufficient information to determine coverage."  
Always cite the exact policy ID and section.  
Question: Does this service require prior authorization based on the retrieved policies?  
Retrieved Policy Chunks:  
{{retrieved\_text}}"*



Ex 2. Sample Prompt for a Claim submission review Workflow

*"You are generating a claims determination rationale.*

*You must base your answer ONLY on the retrieved medical policy text.*

*You may not invent criteria or cite policies that were not retrieved.*

*Always reference policy ID + section number.*

*Task: Explain why the claim was denied, using policy-grounded language.*

*Retrieved Policy Chunks:*

*{{retrieved\_text}}"*

Ex 3. Sample Prompt for Missing documentation review Workflow

*"You are a claims review assistant.*

*Your task is to identify missing documentation required by policy.*

*Use only the retrieved policy extracts provided.*

*If the policy does not explicitly list required documentation, respond with:*

*"Insufficient information in the retrieved policies."*

*Always include citations.*

*Question: What documentation is missing for the submitted claim?*

*Retrieved Policy Chunks:*

*{{retrieved\_text}}"*

## 4. Output Validation

Automated checks enforce:

- Correct citations
- Faithfulness to retrieved text
- No invented rules
- No hallucinations

Tools like RAGAS, custom evaluation pipelines, or SME (Subject Matter Expert) review ensure compliance.

## Why Policy-Grounded RAG Matters in Claims Systems

Claims adjudication is a high-risk process where misinterpretation of policy can lead to:

- Incorrect payments
- Denials without justification
- Non-compliance with CMS, NCQA, URAC, or state regulations
- Time-consuming human reviews
- Rework through resubmission

Policy-Grounded RAG enables safe automation by ensuring AI responses remain:

- ✓ Accurate
- ✓ Audit-friendly
- ✓ Policy-aligned
- ✓ Explainable

**Ex 1. Example: Coverage Rationale**

Based on Policy CP-102: Imaging Services, Section 5.3, MRI without prior authorization is not covered for this indication. The submitted claim lacks prior authorization documentation.

**Ex 2. Example: Denial Explanation**

Policy MP-45: Outpatient Surgery, Section 2.1 requires 24-hour observation notes. These documents were not provided.

In both examples provided, the AI is not reasoning independently—it is quoting actual, validated policy sections.

## Policy Grounded RAG Framework

The Policy-Grounded RAG Implementation Framework consists of six core pillars, each representing a critical domain of capability required for safe enterprise deployment:

1. *Policy Corpus Governance* – Ensuring authoritative, clean, compliant source material
2. *Retrieval Architecture & Metadata Strategy* – High-precision, contextually correct content retrieval
3. *Guard-railed Generation Layer* – Constraining the LLM to eliminate hallucinations
4. *Decision Governance & Validation* – Enforcing policy accuracy and audit readiness
5. *Integration with Claims & Clinical Workflows* – Embedding RAG into real operational systems
6. *Monitoring, Evaluation, and Continuous Improvement* – Sustaining compliance and performance

### 1. Pillar 1 — Policy Corpus Governance

A Policy-Grounded RAG system is only as strong as the data it is allowed to use. Governance ensures safety and regulatory compliance.

#### Key Components

- Authoritative document sources only (medical policies, benefits, SOPs, LOAs)
- Version control and effective-date management
- De-duplication and removal of outdated policy

- Metadata tagging aligned with healthcare operational needs:
  - Plan type (MA, Medicaid, Commercial)
  - State/region
  - LOB
  - Category (imaging, cardiology, surgery)
  - Effective & expiry dates
- SME ownership and approval workflows

**Outcome:**

- ✓ Structured, accurate, regulator-aligned policy corpus.

## 2. Pillar 2 — Retrieval Architecture & Metadata Strategy

Ensures RAG retrieves only the correct, contextually appropriate, and most compliant policy content.

**Key Components**

- Hybrid retrieval (BM25 + embeddings)
- Strict filtering rules (plan, state, LOB, effective date)
- Scoring that prioritizes policy relevance
- Chunking strategy aligned to policy sections
- Low-recall/high-precision retrieval approach for compliance

**Outcome:**

- ✓ AI receives only the correct sections, preventing cross-plan or cross-policy contamination.

### 3. Pillar 3 — Guard railed Generation Layer

Controls how the LLM writes responses and ensures policy faithfulness.

#### Key Guardrails

- System prompts forcing:
  - “Use only retrieved policy text.”
  - “If insufficient information: state uncertainty.”
  - “Always include policy ID + section.”
- No external medical or regulatory knowledge
- Structured output format:

```
{
  "determination": "approve/deny/insufficient_information",
  "rationale": "policy-grounded explanation",
  "citations": ["policy_id:section"]
}
```
- Safety rules preventing:
  - Invented criteria
  - Incorrect citations
  - Non-policy-based reasoning

#### Outcome:

- ✓ Zero-hallucination, audit-ready generation.

#### 4. Pillar 4 — Decision Governance & Validation Layer

Ensures output guarantees policy-accurate, audit-ready explanations that meet CMS, NCQA, and URAC standards..

##### **Validation Components**

- Citation validation engine (checks that cited text exists)
- Faithfulness scoring (does the rationale match retrieved text?)
- SME review sampling (10–20% of outputs)
- Automated blockers during Schema validation for:
  - Missing citations
  - Unsupported reasoning
  - Unretrieved policy references

##### **Outcome:**

- ✓ Outputs meet CMS, NCQA, URAC, and state regulator audit standards.

#### 5. Pillar 5 — Workflow Integration Layer

Deploys Policy-Grounded RAG into actual healthcare operational flows.

##### **Key Integration Areas**

- Claims adjudication (filters, edits, auto-rationales)
- Medical necessity review (clinical criteria extraction)
- Prior authorization workflows
- Appeals and grievances
- Provider communication and call-center copilot tools
- Clinical documentation review

**Integration Methods**

- API-based insertion into adjudication systems
- Human-in-the-loop review layers
- Smart templates for provider-facing content
- RAG outputs embedded into decision support tools

**Outcome:**

- ✓ Consistent, compliant AI across all operational decision points.

**6. Pillar 6 — Monitoring, KPIs & Continuous Improvement**

Ensures the system remains compliant, accurate, and auditable over time.

**Monitoring Metrics (Using RAGAs)**

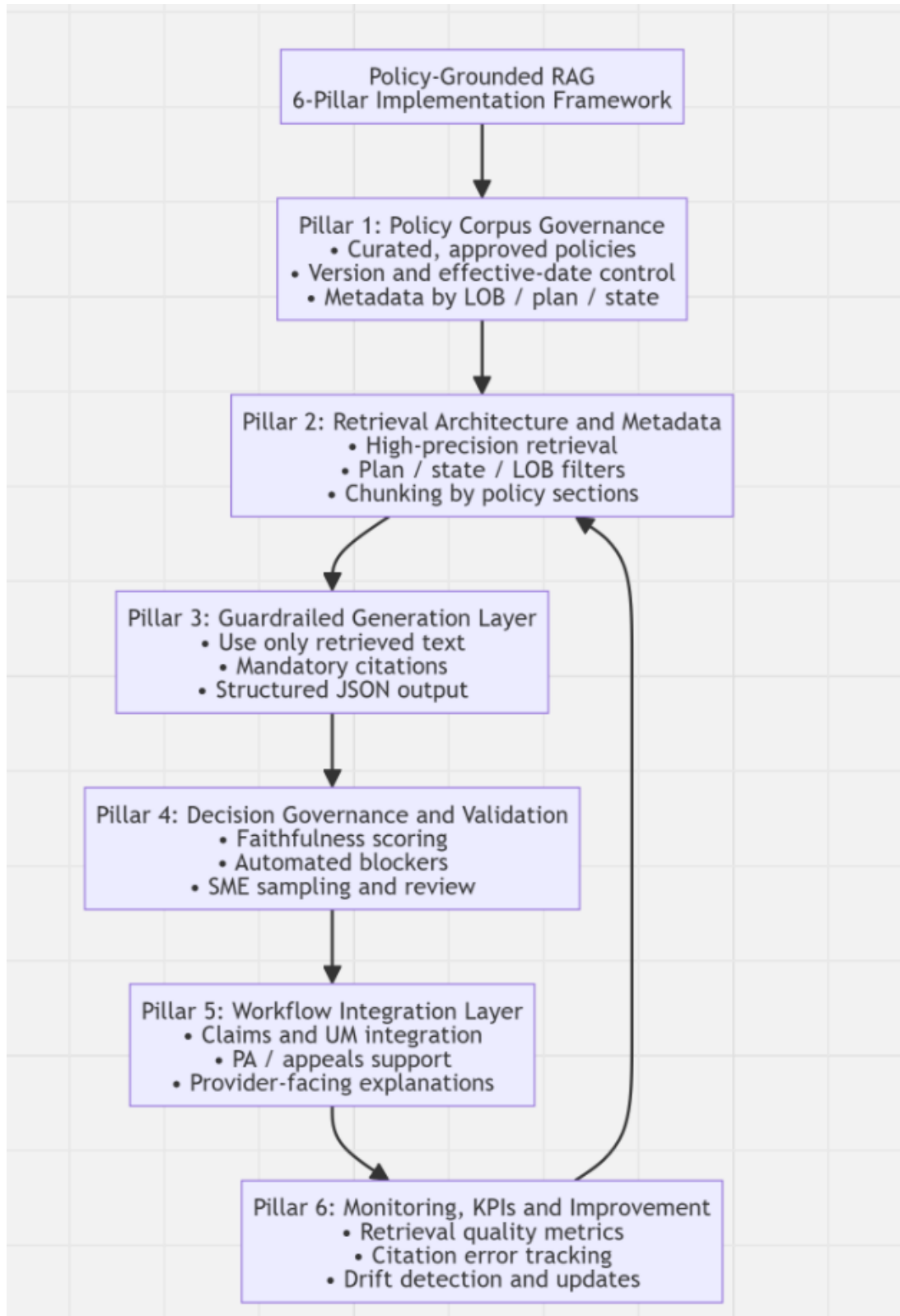
- Faithfulness score –
- Accuracy of policy interpretation
- Turnaround time improvements
- Audit performance (zero-finding targets)
- Error rate in citations
- Retrieval quality metrics

**Continuous Improvement Practices**

- Monthly corpus updates
- Quarterly SME policy audits
- Error-pattern analysis
- Drift detection
- RAG evaluation pipelines (RAGAS, custom metrics)

**Outcome:**

- ✓ A continuously improving, regulator-ready RAG system.





## Technical implementation Blueprint

To implement true Policy-Grounded RAG, organizations must enforce constraints across the entire pipeline.

### 1. Corpus Curation

- Collect only approved documents
- Remove conflicting versions during curation
- Automate pipeline to remove outdated policies and documents
- Gather documents only from trusted sources (directly from regulators such as Government sites or sources)

### 2. Indexing and Metadata Strategy

Index documents with metadata such as:

- Plan type
- LOB
- Effective date
- State
- Category

Chunk documents carefully (policy sections, subsections) to maintain context.

### 3. Strict Retrieval Controls

Use filtering and scoring techniques to ensure:

- High-precision retrieval
- Minimal irrelevant chunks

#### 4. Guardrails and Prompt Engineering

*System prompt:*

"You must answer only using retrieved policy text. If relevant information is missing, respond: 'Insufficient information in the retrieved policy documents.' Always cite policy IDs and sections."

Chunk documents carefully (policy sections, subsections) to maintain context.

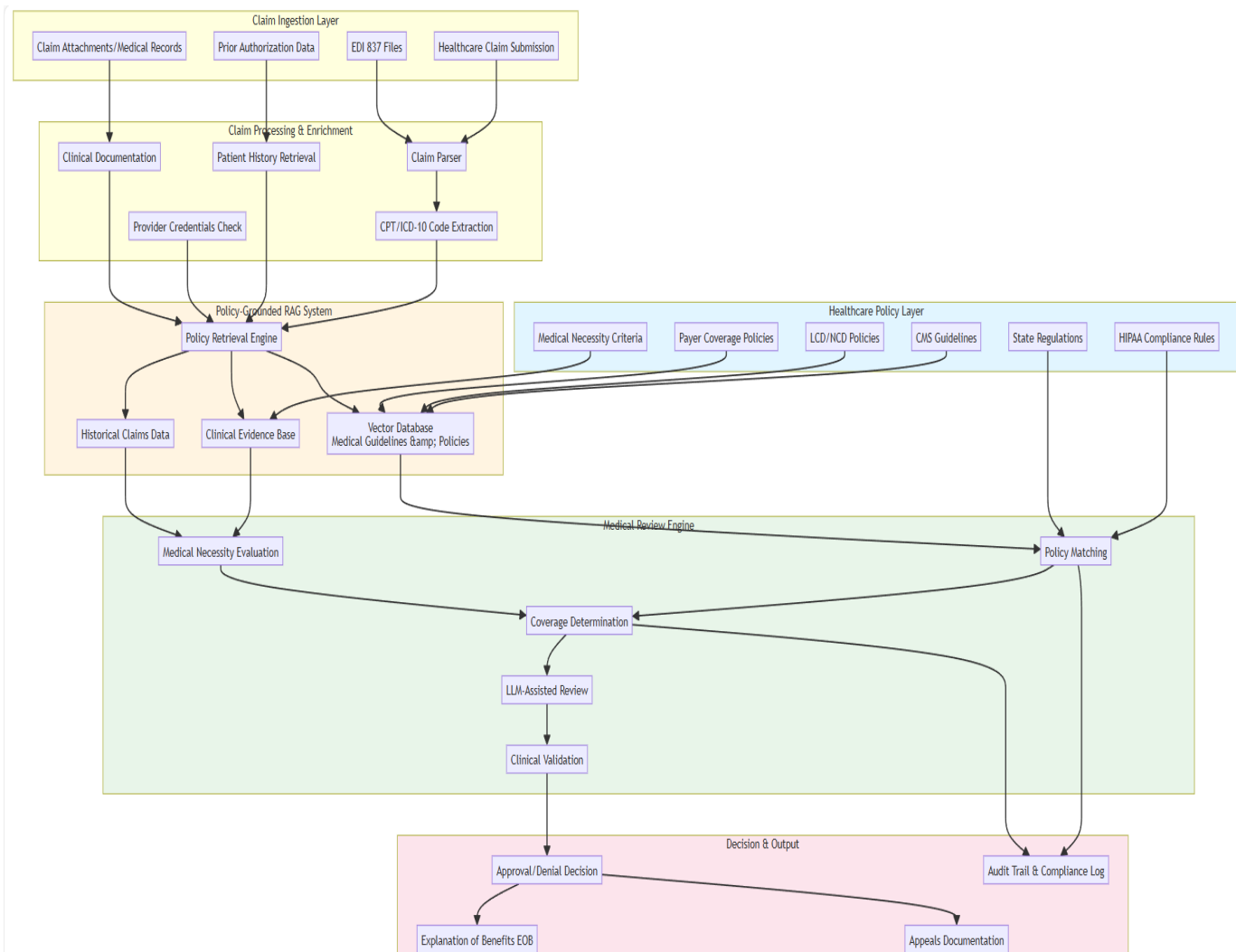
#### 5. Validation and Monitoring

*Implement:*

- Citation schema validation
- Relevance scoring
- Faithfulness checks
- SME sampling and periodic audits

Any answer that does not meet criteria is blocked.

## Sample Policy grounded RAG workflow for an Healthcare claim review



## Benefits of Policy grounded RAG

Policy-Grounded RAG delivers measurable operational, compliance, and financial value especially in claim submission, review, and adjudication workflows. By grounding AI outputs in approved policies, payers and administrators can safely automate decision support while reducing risk.

### 1. Accuracy and Consistency

Policy-Grounded RAG ensures that every explanation, rationale, or coverage review is aligned with the exact wording of medical policies, benefit plans, and SOPs.

- Eliminates inconsistencies across reviewers and teams
- Ensures uniform interpretation of policies
- Removes subjective or variable human judgment in policy lookup

Some of the healthcare Usecase where this applies:

- *Claims Intake*: Accurate identification of required documentation and missing elements
- *Automated Claim Review*: Consistent interpretation of prior authorization, coding, and benefit rules
- *Medical Necessity Review*: Alignment with clinical policy language during determinations

## 2. Compliance, Auditability & Explainability

Grounded answers embed citations to policy IDs, sections, and effective dates. This creates a transparent audit trail. aligned with the exact wording of medical policies, benefit plans, and SOPs.

- Every claim decision can be explained
- Rationale is traceable, reviewable, and defensible
- Supports CMS, NCQA, URAC, and state-level audit requirements
- Gives reference links along with response for quick reference/Manual review etc.

Some of the healthcare Usecase where this applies:

- Decision support: Seamless proof of why a denial/approval was made
- *Post-payment audits*: Why there is variance in claim payment(Underpaid/Overpaid/Write offs)
- *Internal Quality Assurance*: Reviewer coaching and error detection
- *Regulatory audits*: Fast production of policy-backed justification

## 3. Reduced Hallucinations

Because the model cannot rely on general knowledge or invented reasoning, it avoids introducing incorrect, non-policy-based explanations.

- Eliminates fabricated medical rules
- Removes risk of applying policies that do not exist or are out of date
- Strengthens trust in AI-generated review summaries

Some of the healthcare Usecase where this applies:

- *Provider appeals*: Prevents incorrect rationale from escalating issues
- *Initial claim reviews*: Avoids invalid denials created by hallucinated criteria

#### 4. Faster Decisioning

Policy-Grounded RAG pre-generates coverage rationales and missing-information notices grounded in policy text.

- *Significantly reduces time spent manually searching through policies*
- *Speeds up review queues and decreases turnaround times*
- *Improves efficiency in large-volume environments*

Some of the Usecase where this applies in the U.S. healthcare claim lifecycle:

- *Claim edit resolution: Rapid identification of which rule triggered an edit*
- *Medical review triage: Quick drafting of policy-aligned rationales for clinician oversight*
- *Customer service: Instant policy-based responses to provider inquiries*

#### 5. Lower Operational Cost

With automated policy lookup and rationale generation, fewer human resources are needed for repetitive research and explanation tasks.

- *Reduces administrative overhead*
- *Limits expensive escalation pathways*
- *Improves throughput without increasing staff*

Some of the Usecase where this applies in the U.S. healthcare claim lifecycle:

- *Claims Operations: Reduced manual lookup and reduced rework from inconsistent decisions*
- *Provider Relations: Fewer disputes caused by unclear or inconsistent communication*
- *Appeals & Grievances: Lower volume of incorrectly adjudicated claims*

## Traditional Vs Policy-Grounded RAG

Feature / Behavior	Traditional RAG	Policy-Grounded RAG
Data Sources	Mix of open-web, internal, and unverified content	Only curated, approved policy corpus
Risk of Hallucination	Moderate-High	Extremely Low due to strict constraints
Citation Requirements	Optional	Mandatory (policy ID + section)
Compliance Reliability	Inconsistent	Highly aligned with regulatory and policy language
Use in Regulated Domains	Risky without guardrails	Safe for claims, PA, medical review
Explainability	Often unclear; reasoning is based on model inference and may not trace back to source texts	Fully explainable; every statement is backed by specific policy excerpts and citations
Explanation Quality	May generalize or invent rules	Exact quotations and grounded rationales
Suitable for Claim Decisions	Limited	Designed specifically for policy-driven decisions
Audit Readiness	Weak	Strong—traceable and defensible outputs
SME Oversight Required	High	Reduced burden due to accurate grounding

## Use Cases Beyond Claims

Use Cases Beyond Claims Policy-Grounded RAG applies to:

- *Prior authorization automation*
- *Medical necessity review*
- *Customer service co-pilots*
- *Compliance reporting*
- *SOP interpretation and enforcement*

Any workflow requiring policy-based logic can highly benefit from Policy Grounded RAG

## Conclusion

Policy-Grounded RAG is essential for safe and compliant GenAI deployment in regulated domains. By binding AI reasoning to authoritative policy documentation, organizations unlock automation while maintaining trust, transparency, and compliance.

It is the foundation for:

- *Claims automation*
- *Clinical review support*
- *Prior authorization assistance*
- *Provider and member communication*

Organizations adopting this approach gain measurable improvements in accuracy, throughput, and governance