# Mastering RAG

## Bias in Retrieval-Augmented Models

Bias Types and Mitigation Strategies in RAG Systems

**Apply multi-perspective knowledge bases**

4

Enhances retrieval fairness by integrating varied viewpoints.

**Diversify training data sources**

1

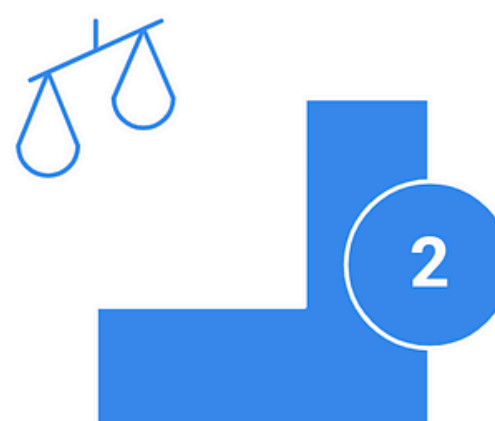Reduces model bias by incorporating diverse perspectives.

**Use counterfactual queries for bias analysis**

3

Identifies user query bias through alternative phrasing.

**Implement fairness constraints in ranking algorithms**

2

Ensures balanced representation in algorithmic decision-making processes.

There are two main layers where bias can enter:

# Bias in the Retrieval Process

- The retriever (often something like a dense retriever or BM25-based system) ranks documents based on relevance.

- But relevance is subjective. It depends on training data, scoring algorithms, and user intent.

- If your retriever is trained on biased clickthrough data (e.g. people tend to click sensational headlines), your model might favor biased or sensational content.

# Bias in the Generation Stage

Even if the retriever pulls neutral or diverse content, the generator (usually an LLM) might:

- Cherry-pick facts to support a certain narrative

- Paraphrase in a way that introduces subtle framing bias

- Omit contradictory viewpoints

# Sources of Bias

- Bias doesn't just come from one place. It's baked into every layer of a RAG system:

| Component | Potential Bias Source |
|---|---|
| Retriever | Biased index data, skewed training queries, over-represented domains |
| Corpus | Inherent societal bias in documents (e.g., underrepresentation of minority voices) |
| Language Model | Pretraining bias from web data, fine-tuning objectives |
| Query Formulation | User framing, lack of context |

# How Bias Manifests in RAG

Let's make it real. Here are some ways bias shows up in practice:

**Overrepresentation Bias:** If your corpus is 80% U.S.-centric, the answers will be too. A RAG model answering health questions might focus on American medical standards, ignoring others.

**Framing Bias**: Even with neutral documents, the language model might generate outputs like:

- "Experts agree that…" (when actually it's debated)
- "Clearly, the best option is…" (even if the evidence is mixed)

**Underexploration Bias**

Retrievers often return the top 5–10 documents. If diverse viewpoints are ranked lower, they may never influence the final answer.

# Why It's Tricky to Fix

Mitigating bias in RAG isn't straightforward because:

- You now have two models to audit: the retriever and the generator.

- Bias in one component can amplify or cancel the other.

- Interventions like re-ranking for fairness may hurt relevance or precision.

- Even defining "fair" retrieval is domain-dependent (e.g., politics vs. science vs. health).

# Possible Mitigations

Here's what researchers and practitioners are exploring:

✅ **Corpus Curation**

- Use balanced datasets with intentional diversity
- Deduplicate sensational or untrustworthy sources

✅ **Retrieval Calibration**

- Incorporate fairness-aware ranking metrics
- Add diversity-promoting loss functions (e.g. Maximal Marginal Relevance)

✅ **Generation Controls**

- Prompt engineering to ask for balanced or multi-perspective answers
- Use classifiers to detect and rewrite biased generations

✅ **Transparency Tools**

- Show retrieved documents alongside answers
- Highlight uncertainty or conflicting evidence