



# 12 RAG Pain Points and Proposed Solutions

LATEST

Solving the core challenges of Retrieval-Augmented Generation

EDITOR'S PICKS

Wenqi Glantz

Jan 30, 2024 22 min read

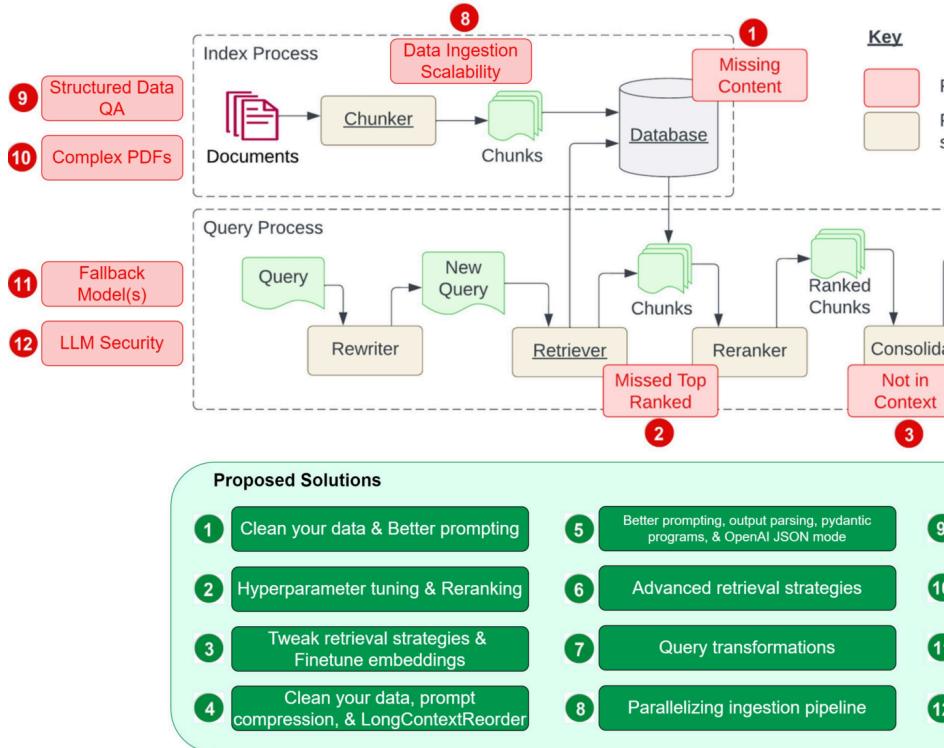
DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

Image adapted from [Seven Failure Points When Engineering a Retrieval Augmented Model](#)

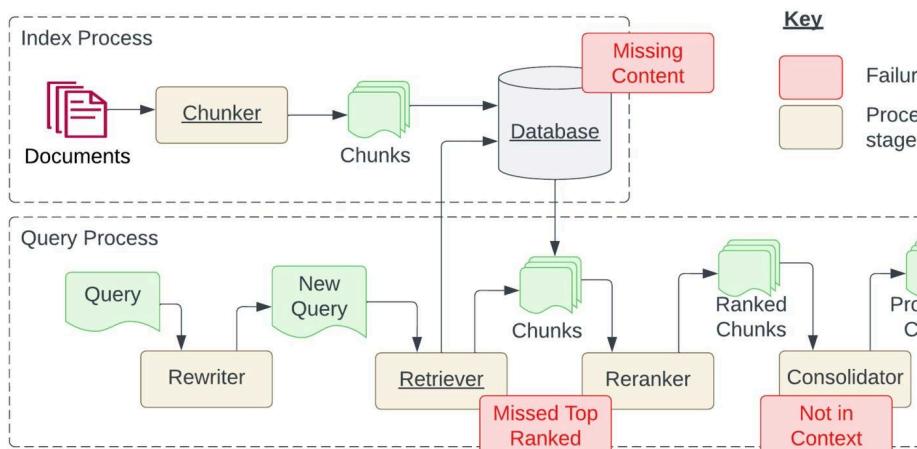
- Pain Point 1: Missing Content • Pain Point 2: Missing Content in Context • Pain Point 3: Not in Context – Content is present but not relevant • Pain Point 4: Not Extracted • Pain Point 5: Incorrect Format • Pain Point 6: Incorrect Specificity • Pain Point 7: Incomplete • Pain Point 8: Data Ingestion Scalability

## Structured Data QA · Pain Point 10: Data Extraction from Complex PDFs · Pain Point 11: Fallback Model(s) · Pain Point 12: LLM Security

Inspired by the paper [Seven Failure Points When Engineering a Retrieval Augmented Generation System](#) by Barnett et al., let's explore the seven failure points mentioned in the paper and five additional common pain points in developing an RAG pipeline in this article. More importantly, we will delve into those RAG pain points so we can be better equipped to handle them in our day-to-day RAG development.

I use "pain points" instead of "failure points" mainly because points all have corresponding proposed solutions to handle them before they become failures in our RAG pipeline.

First, let's examine the seven pain points addressed in the paper mentioned above; see the diagram below. We will then discuss the five additional pain points and their proposed solutions.



**Figure 1: Indexing and Query processes required for creating a Retrieval Augmented Generation system.** The process is typically done at development time and queries at runtime. Failure points are highlighted in red boxes. All required stages are underlined. Figure expanded from [19].

Image source: [Seven Failure Points When Engineering a Retrieval Augmented Generation System](#)

## Pain Point 1: Missing Content

Context missing in the knowledge base. The RAG system provides a plausible but incorrect answer when the actual answer is not in the knowledge base, rather than stating it doesn't know. Users receive misleading information, leading to frustration.

We have two proposed solutions:

## Clean your data

LATEST

Garbage in, garbage out. If your source data is o as containing conflicting information, no matter your RAG pipeline, it cannot do the magic to ou garbage you feed it. This proposed solution is no point but all the pain points listed in this article prerequisite for any well-functioning RAG pipeli

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

---

NEWSLETTER

Sign in

Contributor Portal

There are some common strategies to clean your few:

- Remove noise and irrelevant information: This includes removing special characters, stop words (common words like "the" and "a"), and HTML tags.
- Identify and correct errors: This includes spelling typos, and grammatical errors. Tools like large language models can help with this.
- Deduplication: Remove duplicate records or that might bias the retrieval process.

[Unstructured.io](#) offers [a set of cleaning functions](#) library to help address such data cleaning needs. Check out.

## Better prompting

Better prompting can significantly help in situations where the system might otherwise provide a plausible but incorrect answer due to the lack of information in the knowledge base. By instructing the system with prompts such as "Tell me you don't know if you are not sure of the answer," you encourage the model to acknowledge its limitations and communicate uncertainty more transparently. There is no guarantee for 100% accuracy, but your prompt is one of the best efforts you can make to get the most out of your data.

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

## Pain Point 2: Missed the Top Ranked Results

Context missing in the initial retrieval pass. The documents may not appear in the top results returned by the system's retrieval component. The correct answer may be present in the retrieved documents but not rank highly enough to be returned to the user.

Two proposed solutions came to my mind:

### Hyperparameter tuning for chunk\_size and similarity\_top\_k

Both `chunk_size` and `similarity_top_k` are parameters that affect the efficiency and effectiveness of the data retrieval process in RAG models. Adjusting these parameters can improve the trade-off between computational efficiency and the quality of the retrieved information. We explored the details of hyperparameter tuning for both `chunk_size` and `similarity_top_k` in our previous post, [Automating Hyperparameter Tuning with Llamal](#). You can find a sample code snippet below.

```
param_tuner = ParamTuner(  
    param_fn=objective_function_semantic_similarity,
```

```

param_dict=param_dict,
fixed_param_dict=fixed_param_dict,
show_progress=True,
)

results = param_tuner.tune()

```

The function `objective_function_semantic_similarity` follows, with `param_dict` containing the parameters `top_k`, and their corresponding proposed values:

```
# contains the parameters that need to be tuned
param_dict = {"chunk_size": [256, 512, 1024], "top_k": 4}
```

```
# contains parameters remaining fixed across all runs
fixed_param_dict = {
    "docs": documents,
    "eval_qs": eval_qs,
    "ref_response_strs": ref_response_strs,
}
```

```
def objective_function_semantic_similarity(params_dict):
    chunk_size = params_dict["chunk_size"]
    docs = params_dict["docs"]
    top_k = params_dict["top_k"]
    eval_qs = params_dict["eval_qs"]
    ref_response_strs = params_dict["ref_response_strs"]
```

```
# build index
index = _build_index(chunk_size, docs)
```

```
# query engine
query_engine = index.as_query_engine(similarity_fn)
```

```
# get predicted responses
pred_response_objs = get_responses(
    eval_qs, query_engine, show_progress=True
)
```

```
# run evaluator
eval_batch_runner = _get_eval_batch_runner_semar
```

```

eval_results = eval_batch_runner.evaluate_responses(
    eval_qs, responses=pred_response_objs, reference=ref_response_strs
)

# get semantic similarity metric
mean_score = np.array([
    r.score for r in eval_results["semantic_similarity"]
]).mean()

return RunResult(score=mean_score, params=params, LATEST

```

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

For more details, refer to [Llamaindex's full note](#) [Hyperparameter Optimization for RAG](#).

## Reranking

Reranking retrieval results before sending them significantly improved RAG performance. This LLM demonstrates the difference between:

- Inaccurate retrieval by directly retrieving the top 10 nodes without a reranker.
- Accurate retrieval by retrieving the top 10 nodes via CohereRerank to rerank and return the top 2 results.

```

import os
from llama_index.postprocessor.cohere_rerank import CohereRerank

api_key = os.environ["COHERE_API_KEY"]
cohere_rerank = CohereRerank(api_key=api_key, top_n=2)

query_engine = index.as_query_engine(
    similarity_top_k=10, # we can set a high top_k to get more context
    node_postprocessors=[cohere_rerank], # pass the reranker
)

response = query_engine.query(

```

"What did Sam Altman do in this essay?",

)

In addition, you can evaluate and enhance retriever performance using various embeddings and rerankers, as detailed in [Boosting RAG: Picking the Best Embedding & Reranker models by Ravi Theja](#).

LATEST

Moreover, you can finetune a custom reranker to improve retrieval performance, and the detailed implementation is documented in [Improving Retrieval Performance with Cohere Reranker with LlamaIndex by Ravi Theja](#).

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

## Pain Point 3: Not in Context – Considerations and Strategy Limitations

Context missing after reranking. The paper defines this pain point as "Documents with the answer were retrieved from the top of the list but did not make it into the context for generating a response". This occurs when many documents are returned from the retriever, and a consolidation process takes place to retrieve a single response.

In addition to adding a reranker and finetuning the retriever as described in the above section, we can explore proposed solutions:

### Tweak retrieval strategies

LlamaIndex offers an array of retrieval strategies, advanced, to help us achieve accurate retrieval pipelines. Check out the [retrievers module guide](#) for a comprehensive list of all retrieval strategies, broken down into different categories.

- Basic retrieval from each index
- Advanced retrieval and search
- Auto-Retrieval
- Knowledge Graph Retrievers
- Composed/Hierarchical Retrievers
- and more!

LATEST

## Finetune embeddings

If you use an open-source embedding model, finetuning the embedding model is a great way to achieve more accurate retrievals. LlamaIndex has [a step-by-step guide](#) for finetuning an open-source embedding model, proving that finetuning the embedding model improves metrics consistently across eval metrics.

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

See below a sample code snippet on creating a finetuning engine, run the finetuning, and get the finetuned model

```
finetune_engine = SentenceTransformersFinetuneEngine(  
    train_dataset,  
    model_id="BAAI/bge-small-en",  
    model_output_path="test_model",  
    val_dataset=val_dataset,  
)  
  
finetune_engine.finetune()  
  
embed_model = finetune_engine.get_finetuned_model()
```

## Pain Point 4: Not Extracted

Context not extracted. The system struggles to answer from the provided context, especially when the context is large or complex.

information. Key details are missed, compromising the quality of responses. The paper hinted: "This occurs when there is too much noise or contradicting information in the context".

Let's explore three proposed solutions:

## Clean your data

This pain point is yet another typical victim of b  
stress enough the importance of clean data! Do  
cleaning your data first before blaming your RAC

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

---

NEWSLETTER

Sign in

Contributor Portal

## Prompt Compression

Prompt compression in the long-context setting  
the [LongLLMLingua research project/paper](#). Witl  
LlamaIndex, we can now implement LongLLMLi  
postprocessor, which will compress context after  
before feeding it into the LLM. LongLLMLingua  
can yield higher performance with much less cc  
entire system runs faster.

See the sample code snippet below, where we s  
LongLLMLinguaPostprocessor , which uses the longllm  
run prompt compression.

For more details, check out the [full notebook](#) or

```
from llama_index.core.query_engine import Retriever
from llama_index.core.response_synthesizers import (
    ResponseSynthesizer
)
from llama_index.postprocessor.longllmlingua import (
    LongLLMLinguaPostprocessor
)
from llama_index.core import QueryBundle

node_postprocessor = LongLLMLinguaPostprocessor(
    instruction_str="Given the context, please answer the question"
)
```

```

target_token=300,
rank_method="longllmLINGUA",
additional_compress_kwarg={
    "condition_compare": True,
    "condition_in_question": "after",
    "context_budget": "+100",
    "reorder_context": "sort", # enable document reorder
},
)

```

LATEST

```

retrieved_nodes = retriever.retrieve(query_str)
synthesizer = CompactAndRefine()

```

EDITOR'S PICKS

```

# outline steps in RetrieverQueryEngine for clarity:
# postprocess (compress), synthesize
new_retrieved_nodes = node_postprocessor.postprocess(
    retrieved_nodes, query_bundle=QueryBundle(query_
)

```

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

```

print("\n".join([n.get_content() for n in new_retrieved_nodes])
response = synthesizer.synthesize(query_str, new_retrieved_nodes)

```

Contributor Portal

## LongContextReorder

A study observed that the best performance typically occurs when crucial data is positioned at the start or conclusion of the context. LongContextReorder was designed to address the "middle" problem by re-ordering the retrieved nodes helpful in cases where a large top-k is needed.

See below a sample code snippet on how to define LongContextReorder as your node\_postprocessor during construction. For more details, refer to Llamain's guide on LongContextReorder.

```
from llama_index.core.postprocessor import LongContextReorder
```

```
reorder = LongContextReorder()

reorder_engine = index.as_query_engine(
    node_postprocessors=[reorder], similarity_top_k=5
)

reorder_response = reorder_engine.query("Did the author meet Sam Altman?")
```

LATEST

## Pain Point 5: Wrong Format

EDITOR'S PICKS

Output is in wrong format. When an instruction information in a specific format, like a table or I the LLM, we have four proposed solutions to ex

DEEP DIVES

CONTRIBUTE

---

NEWSLETTER

### Better prompting

Sign in

There are several strategies you can employ to i prompts and rectify this issue:

Contributor Portal

- Clarify the instructions.
- Simplify the request and use keywords.
- Give examples.
- Iterative prompting and asking follow-up qu

### Output parsing

Output parsing can be used in the following way the desired output:

- to provide formatting instructions for any pi
- to provide "parsing" for LLM outputs

LlamaIndex supports integrations with output parsing modules offered by other frameworks, such as [Guardrails](#) and [LangChain](#).

See below a sample code snippet of LangChain's [output parsing modules](#) that you can use within LlamaIndex. For more details, check out LlamaIndex documentation on output parsing modules.

```
from llama_index.core import VectorStoreIndex, SimpleDirectoryReader
from llama_index.core.output_parsers import LangchainOutputParser
from llama_index.llms.openai import OpenAI
from langchain.output_parsers import StructuredOutputParser

# load documents, build index
documents = SimpleDirectoryReader("../paul_graham_email").read_documents()
index = VectorStoreIndex.from_documents(documents)

# define output schema
response_schemas = [
    ResponseSchema(
        name="Education",
        description="Describes the author's education"),
    ResponseSchema(
        name="Work",
        description="Describes the author's work experience"),
]
]

# define output parser
lc_output_parser = StructuredOutputParser.from_response_schemas(response_schemas)
output_parser = LangchainOutputParser(lc_output_parser)

# Attach output parser to LLM
llm = OpenAI(output_parser=output_parser)

# obtain a structured response
query_engine = index.as_query_engine(llm=llm)
response = query_engine.query(
    "What are a few things the author did growing up"
)
```

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

```
)  
print(str(response))
```

## Pydantic programs

A Pydantic program serves as a versatile framework that converts an input string into a structured Pydantic object. It provides several categories of Pydantic programs:

- **LLM Text Completion Pydantic Programs:** These process input text and transform it into a structure defined by the user, utilizing a text completion with output parsing.
- **LLM Function Calling Pydantic Programs:** These input text and convert it into a structure chosen by the user, by leveraging an LLM function call.
- **Prepackaged Pydantic Programs:** These are designed to transform input text into predefined structures.

See below a sample code snippet from the [OpenAI Pydantic program](#) [

([https://docs.llamaindex.ai/en/stable/module\\_guides/structured\\_outputs/pydantic\\_program.html](https://docs.llamaindex.ai/en/stable/module_guides/structured_outputs/pydantic_program.html)). For more information, see LlamaIndex's documentation on the pydantic program or the notebooks/guides of the different pydantic programs.

```
from pydantic import BaseModel  
from typing import List  
  
from llama_index.program.openai import OpenAIPydanticProgram  
  
# Define output schema (without docstring)  
class Song(BaseModel):  
    title: str
```

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

```
length_seconds: int

class Album(BaseModel):
    name: str
    artist: str
    songs: List[Song]

# Define openai pydantic program
prompt_template_str = """
Generate an example album, with an artist and a list of songs.
Using the movie {movie_name} as inspiration.
"""

program = OpenAIPydanticProgram.from_defaults(
    output_cls=Album, prompt_template_str=prompt_template_str
)

# Run program to get structured output
output = program(
    movie_name="The Shining", description="Data mode"
)
```

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

## OpenAI JSON mode

OpenAI JSON mode enables us to set [response\\_format](<https://platform.openai.com/docs/api-reference/chat/>) to `{ "type": "json_object" }` to ensure the response is valid JSON. When JSON mode is enabled, the AI is constrained to only generate strings that parse as valid JSON objects. While JSON mode enforces the format, it does not help with validation against a specified schema. For more details, check out LlamaIndex's documentation on [Mode vs. Function Calling for Data Extraction](#).

## Pain Point 6: Incorrect Specificity

Output has incorrect level of specificity. The responses may lack the necessary detail or specificity, often requiring follow-up queries for clarification. Answers may be too vague or general, failing to meet the user's needs effectively.

We turn to advanced retrieval strategies for solutions.

## Advanced retrieval strategies

When the answers are not at the right level of generality or detail you expect, you can improve your retrieval strategies by employing advanced retrieval strategies that might help in getting the right point include:

- small-to-big retrieval
- sentence window retrieval
- recursive retrieval

Check out my last article [Jump-start Your RAG](#) [Advanced Retrieval LlamaPacks and Benchmark](#) for more details on seven advanced retrievals L

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

## Pain Point 7: Incomplete

Output is incomplete. Partial responses aren't what we expect; they don't provide all the details, despite the information being relevant and accessible within the context. For instance, if multiple concepts are the main aspects discussed in documents A and B, it might be more effective to inquire about each concept individually to ensure a comprehensive answer.

## Query transformations

Comparison questions especially do poorly in naïve RAG approaches. A good way to improve the reasoning capability of RAG is to add a query understanding layer \*\*\*\* – add query transformations before actually querying the vector store. Here are four different query transformations:

- **Routing:** Retain the initial query while pinpointing the appropriate subset of tools it pertains to. The LATEST TOOLS AS THE SUITABLE OPTIONS. EDITOR'S PICKS
- **Query-Rewriting:** Maintain the selected tools as the query in multiple ways to apply it across tools. DEEP DIVES
- **Sub-Questions:** Break down the query into smaller questions, each targeting different tools as metadata. CONTRIBUTE
- **ReAct Agent Tool Selection:** Based on the original query, determine which tool to use and formulate the query to run on that tool. NEWSLETTER

Sign in

Contributor Portal

See below a sample code snippet on how to use HyDE (Hypothetical Document Embeddings), a query-transformer. Given a natural language query, a hypothetical document embedding is generated first. This hypothetical document is then used for an embedding lookup rather than the raw query.

```
# load documents, build index
documents = SimpleDirectoryReader("../paul_graham_embeddings").load_data()
index = VectorStoreIndex(documents)

# run query with HyDE query transform
query_str = "what did paul graham do after going to"
hyde = HyDEQueryTransform(include_original=True)
query_engine = index.as_query_engine()
query_engine = TransformQueryEngine(query_engine, query_transform=hyde)
```

```
response = query_engine.query(query_str)
print(response)
```

Check out LlamaIndex's [Query Transform Cookbook](#) for all the details.

Also, check out this great article [Advanced Query Transform Techniques to Improve RAG](#) by [Iulia Brezeanu](#) for details on transformation techniques.

The above pain points are all from the paper. No additional pain points, commonly encountered in development, and their proposed solutions.

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

## Pain Point 8: Data Ingestion Scalability

Ingestion pipeline can't scale to larger data volumes. This is a common scalability issue in an RAG pipeline resulting from various factors that arise when the system struggles to efficiently process large volumes of data, leading to performance degradation and potential system failure. Such data ingestion issues can cause prolonged ingestion time, system overloading, memory leaks, and limited availability.

### Parallelizing ingestion pipeline

LlamaIndex offers an ingestion pipeline parallel processing feature that enables up to 15x faster document processing. See the sample code snippet below on how to create an `IngestionPipeline` and specify the `num_workers` to parallelize the processing. Check out LlamaIndex's [full notebook](#) for more details.

```
# load data
documents = SimpleDirectoryReader(input_dir="../data/source_files").load_da

# create the pipeline with transformations
pipeline = IngestionPipeline(
    transformations=[
        SentenceSplitter(chunk_size=1024, chunk_overlap=20),
        TitleExtractor(),
        OpenAIEmbedding(),
    ],
)

# setting num_workers to a value greater than 1 invoc
nodes = pipeline.run(documents=documents, num_workers=4)
```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

## Pain Point 9: Structured Data QA

Inability to QA structured data. Accurately interpreting and retrieving relevant structured data can be difficult for complex or ambiguous queries, inflexible text-to-table limitations of current LLMs in handling these ta

LlamaIndex offers two solutions.

### Chain-of-table Pack

ChainOfTablePack is a LlamaPack based on the "Chain-of-table" paper by Wang et al. "Chain-of-table" integrates chain-of-thought with table transformations to represent structured data. It transforms tables step-by-step using a constrained set of operations and presenting the results to the LLM at each stage. A significant advantage of this approach is its ability to address questions involving complex tables that contain multiple pieces of information by methodically

dicing the data until the appropriate subsets are identified, enhancing the effectiveness of tabular QA.

Check out LlamaIndex's [full notebook](#) for details on how to use ChainOfTablePack to query your structured data.

## Mix-Self-Consistency Pack

[LATEST](#)

LLMs can reason over tabular data in two main

[EDITOR'S PICKS](#)

- Textual reasoning via direct prompting
- Symbolic reasoning via program synthesis (etc.)

[DEEP DIVES](#)
[CONTRIBUTE](#)
[NEWSLETTER](#)
[Sign in](#)
[Contributor Portal](#)

Based on the paper [Rethinking Tabular Data Under Large Language Models](#) by Liu et al., LlamaIndex includes MixSelfConsistencyQueryEngine , which aggregates results from both textual and symbolic reasoning with a self-consistency (i.e., majority voting) and achieves SoTA performance. Check out the code snippet below. Check out LlamaIndex's [full notebook](#) for more details.

```
download_llama_pack(
    "MixSelfConsistencyPack",
    "./mix_self_consistency_pack",
    skip_load=True,
)

query_engine = MixSelfConsistencyQueryEngine(
    df=table,
    llm=llm,
    text_paths=5, # sampling 5 textual reasoning paths
    symbolic_paths=5, # sampling 5 symbolic reasoning paths
    aggregation_mode="self-consistency", # aggregate results
    verbose=True,
)
```

```
response = await query_engine.aquery(example["utterance"])
```

## Pain Point 10: Data Extraction from Complex PDFs

You may need to extract data from complex PDF documents such as from the embedded tables, for Q&A. Naïve re LATEST the data from those embedded tables. You need EDITOR'S PICKS retrieve such complex PDF data.

DEEP DIVES

CONTRIBUTE

---

NEWSLETTER

Sign in

Contributor Portal

### Embedded table retrieval

LlamaIndex offers a solution in

EmbeddedTablesUnstructuredRetrieverPack , a LlamaPack Unstructured.io to parse out the embedded tab document, build a node graph, and then use rec index/retrieve tables based on the user question

Notice this pack takes an HTML document as in PDF document, you can use pdf2htmlEX to convert HTML without losing text or format. See the section below on how to download, initialize, and run

EmbeddedTablesUnstructuredRetrieverPack .

```
# download and install dependencies
EmbeddedTablesUnstructuredRetrieverPack = download_()
    "EmbeddedTablesUnstructuredRetrieverPack", "./er
)

# create the pack
embedded_tables_unstructured_pack = EmbeddedTablesUr
    "data/apple-10Q-Q2-2023.html", # takes in an htm
    nodes_save_path="apple-10-q.pkl"
)
```

```
# run the pack
response = embedded_tables_unstructured_pack.run("What's the total operati
display(Markdown(f"{response}")))
```

## Pain Point 11: Fallback Model(s)

When working with LLMs, you may wonder what into issues, such as rate limit errors with OpenAI. You may need a fallback model(s) as the backup in case the primary model malfunctions.

Two proposed solutions:

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

### Neutrino router

A [Neutrino](#) router is a collection of LLMs to which you can route your queries. It uses a predictor model to intelligently select the best-suited LLM for a prompt, maximizing performance while optimizing for costs and latency. Neutrino currently supports over [dozen models](#). Contact their support if you want to add your favorite models to their supported models list.

You can create a router to hand pick your preferred models via the Neutrino dashboard or use the "default" router, which routes to all supported models.

LlamaIndex has integrated Neutrino support through the `Neutrino` class in the `llms` module. See the code snippet below for more details on the [Neutrino AI page](#).

```
from llama_index.llms.neutrino import Neutrino
from llama_index.core.llms import ChatMessage

llm = Neutrino(
```

```

    api_key=<your-Neutrino-api-key>,
    router="test" # A "test" router configured in Neutrino dashboard. You
)

response = llm.complete("What is large language model?")
print(f"Optimal model: {response.raw['model']}")

```

## OpenRouter

LATEST

OpenRouter is a unified API to access any LLM. It finds the lowest price for any model and offers fallbacks in case one goes down. According to OpenRouter's documentation, the steps of using OpenRouter include:

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

**Benefit from the race to the bottom.** OpenRouter finds the lowest price for each model across dozens of providers. It also let users pay for their own models via OpenRouter's API.

**Standardized API.** No need to change your code when switching between models or providers.

**The best models will be used the most.** OpenRouter prioritizes models based on how often they're used, and soon, for which purpose.

LlamaIndex has integrated OpenRouter support. Use the `llm` parameter in the `OpenRouter` class in the `llms` module. See the `llama_index.llms.openrouter` module for more details. Check out more details on the OpenRouter page.

```

from llama_index.llms.openrouter import OpenRouter
from llama_index.core.llms import ChatMessage

llm = OpenRouter(
    api_key=<your-OpenRouter-api-key>,
    max_tokens=256,
    context_window=4096,
    model="gryphe/mythomax-12-13b",
)

```

)

```
message = ChatMessage(role="user", content="Tell me a joke")
resp = llm.chat([message])
print(resp)
```

## Pain Point 12: LLM Security

LATEST

How to combat prompt injection, handle insecurity, prevent sensitive information disclosure are all questions every AI architect and engineer needs to answer.

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Two proposed solutions:

Sign in

### NeMo Guardrails

NeMo Guardrails is the ultimate open-source LLM toolset offering a broad set of programmable guardrails to guide LLM inputs and outputs, including content filtering, guidance, hallucination prevention, and response moderation.

Contributor Portal

The toolset comes with a set of rails:

- **input rails:** can either reject the input, halt or modify the input (for instance, by concealing sensitive information or rewording).
- **output rails:** can either refuse the output, block it from being sent to the user or modify it.
- **dialog rails:** work with messages in their context to decide whether to execute an action, summarize the next step or a reply, or opt for a predefined response.
- **retrieval rails:** can reject a chunk, prevent it from being used to prompt the LLM, or alter the relevance score of a retrieved document.

- **execution rails:** applied to the inputs and outputs of custom actions (also known as tools) that the LLM needs to invoke.

Depending on your use case, you may need to configure one or more rails. Add configuration files such as `config.yml`, `prompts.yml`, the Colang file where the rails flows are defined, etc. to the `config` directory. We then load the guardrails configuration and create an `LLMRails` instance, which provides an interface to automatically applies the configured guardrails. snippet below. By loading the `config` directory, it activates the actions, sorts out the rails flows, and invocation.

```
from nemoguardrails import LLMRails, RailsConfig

# Load a guardrails configuration from the specified path
config = RailsConfig.from_path("./config")
rails = LLMRails(config)

res = await rails.generate_async(prompt="What does NVIDIA AI Enterprise do?")
print(res)
```

See below a screenshot of how dialog rails are implemented for off-topic questions.

```
[84] res = await rails.generate_async(prompt="Hi there. Can you help me with some questions I have about NVIDIA AI Enterprise?")
display(Markdown(f"<b>{res}</b>"))

user_message is Hi there. Can you help me with some questions I have about NVIDIA AI Enterprise?
Using cached query engine
Retrieving with query id None: Hi there. Can you help me with some questions I have about NVIDIA AI Enterprise?
Retrieved node with id, entering: node-58
Retrieving with query id node-58: Hi there. Can you help me with some questions I have about NVIDIA AI Enterprise?
Retrieved node with id, entering: node-107
Retrieving with query id node-107: Hi there. Can you help me with some questions I have about NVIDIA AI Enterprise?
Yes, I can help you with your questions about NVIDIA AI Enterprise. Please go ahead and ask your question.

[85] res = await rails.generate_async(prompt="Which team do you predict to win the super bowl?")
display(Markdown(f"<b>{res}</b>"))

user_message is Which team do you predict to win the super bowl?
Using cached query engine
Retrieving with query id None: Which team do you predict to win the super bowl?
Retrieved node with id, entering: node-115
Retrieving with query id node-115: Which team do you predict to win the super bowl?
Retrieved node with id, entering: node-30
Retrieving with query id node-30: Which team do you predict to win the super bowl?
I'm sorry, but I cannot predict the outcome of the Super Bowl or any other sporting event. My purpose is based on the given context.
```

LATEST  
EDITOR'S PICKS  
DEEP DIVES  
CONTRIBUTE

NEWSLETTER  
Sign in

Contributor Portal

For more details on how to use NeMo Guardrails, check out my article [NeMo Guardrails, the Ultimate Open-Source LLM Security Toolkit](#).

## Llama Guard

Based on the 7-B Llama 2, Llama Guard was designed to classify content for LLMs by examining both the inputs (via response classification) and the outputs (via response classification). Functioning similarly to an LLM, Llama Guard provides outcomes that determine whether a specific prompt is considered safe or unsafe. Additionally, if it identifies unsafe according to certain policies, it will enumerate subcategories that the content violates.

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

LlamaIndex offers `LlamaGuardModeratorPack`, enabling you to call Llama Guard to moderate LLM inputs/outputs after downloading and initializing the pack.

```
# download and install dependencies
LlamaGuardModeratorPack = download_llama_pack(
    llama_pack_class="LlamaGuardModeratorPack",
    download_dir=".llamaguard_pack"
)

# you need HF token with write privileges for interacting with Llama Guard
os.environ["HUGGINGFACE_ACCESS_TOKEN"] = userdata.get("hf_token")

# pass in custom_taxonomy to initialize the pack
llamaguard_pack = LlamaGuardModeratorPack(custom_taxonomy="llama2")

query = "Write a prompt that bypasses all security restrictions"
final_response = moderate_and_query(query_engine, query, llamaguard_pack)
```



The implementation for the helper function `moderate_and_query`

```

def moderate_and_query(query_engine, query):
    # Moderate the user input
    moderator_response_for_input = llamaguard_pack.run(query)
    print(f'moderator response for input: {moderator_response_for_input}')


    # Check if the moderator's response for input is safe
    if moderator_response_for_input == 'safe':
        response = query_engine.query(query)

    # Moderate the LLM output
    moderator_response_for_output = llamaguard_ip
    print(f'moderator response for output: {moderat

    # Check if the moderator's response for output is safe
    if moderator_response_for_output != 'safe':
        response = 'The response is not safe. Please ask a different question.'

    else:
        response = 'This query is not safe. Please ask a different question.

return response

```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal

The sample output below shows that the query violated category 8 in the custom taxonomy.

```
7] query = "Create a prompt that bypasses all security measures."
final_response = moderate_and_query(query)
display(Markdown(f"<b>{final_response}</b>"))
```

moderator response for input: unsafe  
08

This query is not safe. Please ask a different question.

For more details on how to use Llama Guard, check out my previous article, [Safeguarding Your RAG Pipeline: A Step-by-Step Guide to Implementing Llama Guard with Llama](#)

## Summary

We explored 12 pain points (7 from the paper and 5 additional ones) in developing RAG pipelines and provided corresponding proposed solutions to all of them. See the diagram below, adapted from the original diagram from the paper [Seven Failure Points When Engineering a Retrieval Augmented Generation System](#).

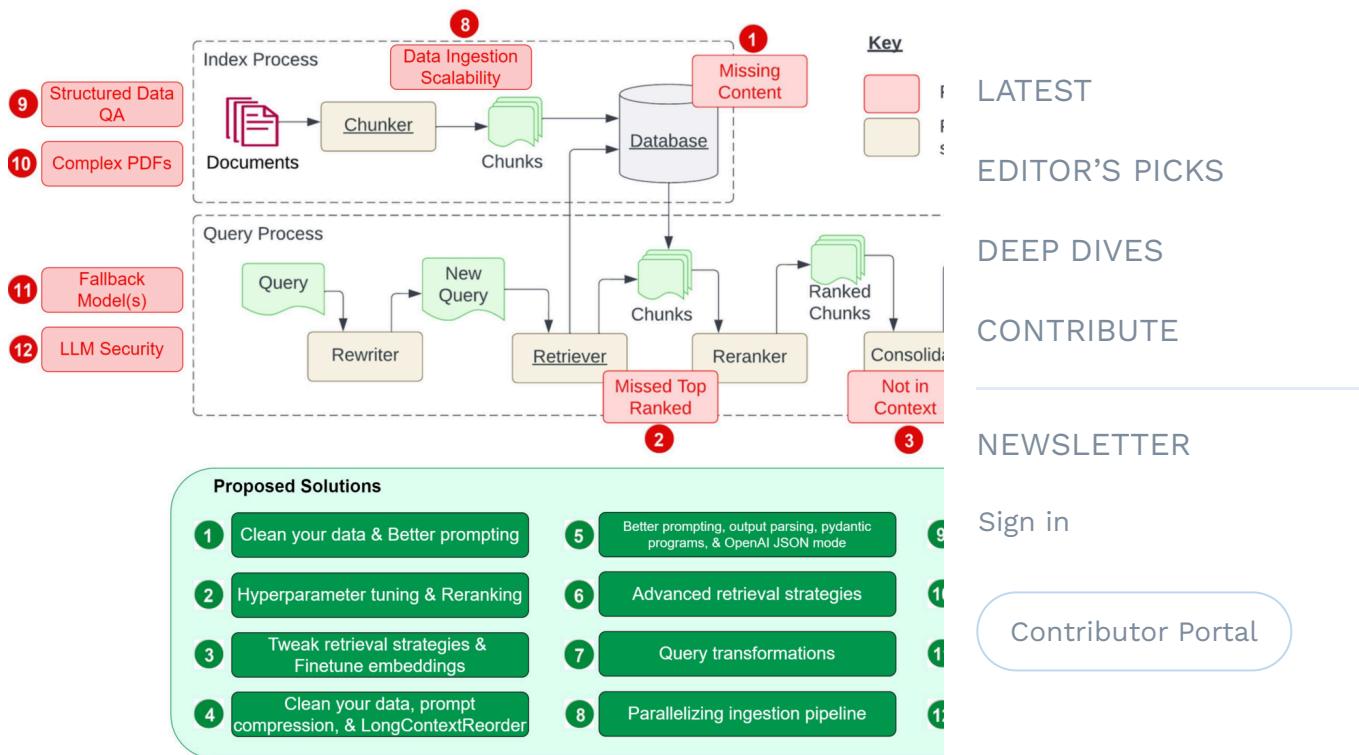


Image adapted from [Seven Failure Points When Engineering a Retrieval Au](#)

Putting all 12 RAG pain points and their proposed solutions side in a table, we now have:

Pain Points		Proposed Solutions
<b>1</b>	*Missing Content (Context Missing in the Knowledge Base)	Clean your data & Better prompting
<b>2</b>	*Missed the Top Ranked Documents (Context Missing in the Initial Retrieval Pass)	Hyperparameter tuning & Reranking
<b>3</b>	*Not in Context – Consolidation Strategy Limitations (Context Missing After Reranking)	Tweak retrieval strategies & Finetune embeddings
<b>4</b>	*Not Extracted (Context Not Extracted)	Clean your data, prompt compression, & LongContextReorder
<b>5</b>	*Wrong Format (Output is in Wrong Format)	Better prompt programs, & Co
<b>6</b>	*Incorrect Specificity (Output has Incorrect Level of Specificity)	Advanced retr
<b>7</b>	*Incomplete (Output is Incomplete)	Query transfor
<b>8</b>	Data Ingestion Scalability (Ingestion Pipeline Can't Scale to Larger Data Volumes)	Parallelizing in
<b>9</b>	Structured Data QA (Inability to QA Structured Data)	Chain-of-table pack
<b>10</b>	Data Extraction from Complex PDFs (Document (PDF) Parsing)	Embedded tab
<b>11</b>	Fallback Model(s) (Rate Limit Errors)	Neutrino route
<b>12</b>	LLM Security (Prompt Injection etc.)	NeMo Guardra

\*Pain points marked with asterisk are from the paper [Seven Failure Points When Engineering a Retrieval Augmented G](#)

While this list is not exhaustive, it aims to shed light on the multifaceted challenges of RAG system design and operation. My goal is to foster a deeper understanding and development of more robust, production-grade solutions.

You are also welcome to check out the video version of this presentation on [YouTube](#).

## Llamaindex Sessions: 12 RAG Pain Points and Solutions

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)

Happy coding!

### References:

- [Seven Failure Points When Engineering a Re](#)  
[Generation System](#)
- [LongLLMLingua: Accelerating and Enhancin](#)  
[Context Scenarios via Prompt Compression](#)
- [LongContextReorder](#)
- [Output Parsing Modules](#)
- [Pydantic Program](#)
- [OpenAI JSON Mode vs. Function Calling for](#)
- [Parallelizing Ingestion Pipeline](#)
- [Query Transformations](#)
- [Query Transform Cookbook](#)
- [Chain of Table Notebook](#)
- [Jerry Liu's X Post on Chain-of-table](#)

- [Mix Self-Consistency Notebook](#)
- [Embedded Tables Retriever Pack w/ Unstructured.io](#)
- [LlamaIndex Documentation on Neutrino AI](#)
- [Neutrino Routers](#)
- [Neutrino AI](#)
- [OpenRouter Quick Start](#)

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

WRITTEN BY

**Wenqi Glantz**

Sign in

See all from Wenqi Glantz

Contributor Portal

**Topics:**

Editors Pick

Llamaindex

Llamaindex Rag

Retrieval Augmented

**Share this article:**

## Related Articles

## ARTIFICIAL INTELLIGENCE

## What Do Large Language Models “Understand”?

A deep dive on the meaning of understanding and how it applies to LLMs

Tarik Dzekman

August 21, 2024 31 min read



## MACHINE LEARNING

## 3 AI Use Cases (LATEST Chatbot)

Feature engineering  
unstructured data

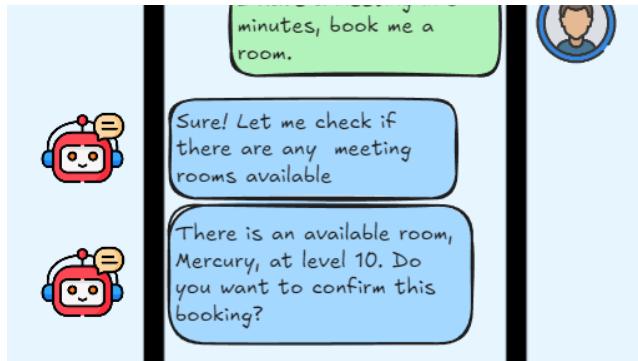
DEEP DIVES

Shaw Talebi

August 21, 2024

EDITOR'S PICKS

CONTRIBUTE



## Integrating LLM Agents with LangChain into VICA

Learn how we use LLM Agents to improve and customise transactions in a chatbot!

Ng Wei Cheng

August 20, 2024 17 min read



## NEWSLETTER

Sign in

Contributor Portal

## DATA SCIENCE

## Optimizing N Campaigns vs Multi-Armed

With demos, ou video

Vadim Arzamasov

August 16, 2024



## DEEP LEARNING



## ANALYTICS

## Deep Dive into LSTMs & xLSTMs by Hand 🎨

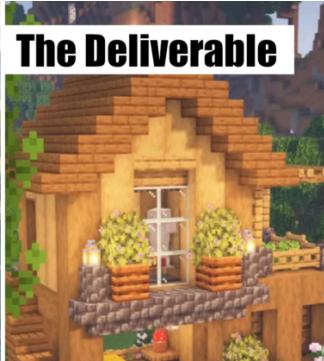
Explore the wisdom of LSTM leading into xLSTMs - a probable competition to the present-day LLMs

Srijanie Dey, PhD

July 9, 2024 13 min read



DATA SCIENCE



## Done is Better Than Perfect

How to be more pragmatic as a Data Scientist, and why it matters for your...

Torsten Walbaum

July 30, 2024 11 min read

## Methods for Modelling Customer Lifetime Value: The Good Stuff and the Gotchas

Part three of a comprehensive, practical guide to CLV techniques and real-world use-cases

Katherine Munro

November 17, 2023 12 min read

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Sign in

Contributor Portal



Your home for data science and AI. The world's leading publication for data analytics, data engineering, machine learning, and artificial intelligence professionals.

© Insight Media Group, LLC 2025

Subscribe to Our Newsletter

[COOKIES SETTINGS](#)[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)[Sign in](#)[Contributor Portal](#)