



12 RAG Pain Points and Proposed Solutions

LATEST

Solving the core challenges of Retrieval-Augmented Ge EDITOR'S PICKS

Wenqi Glantz

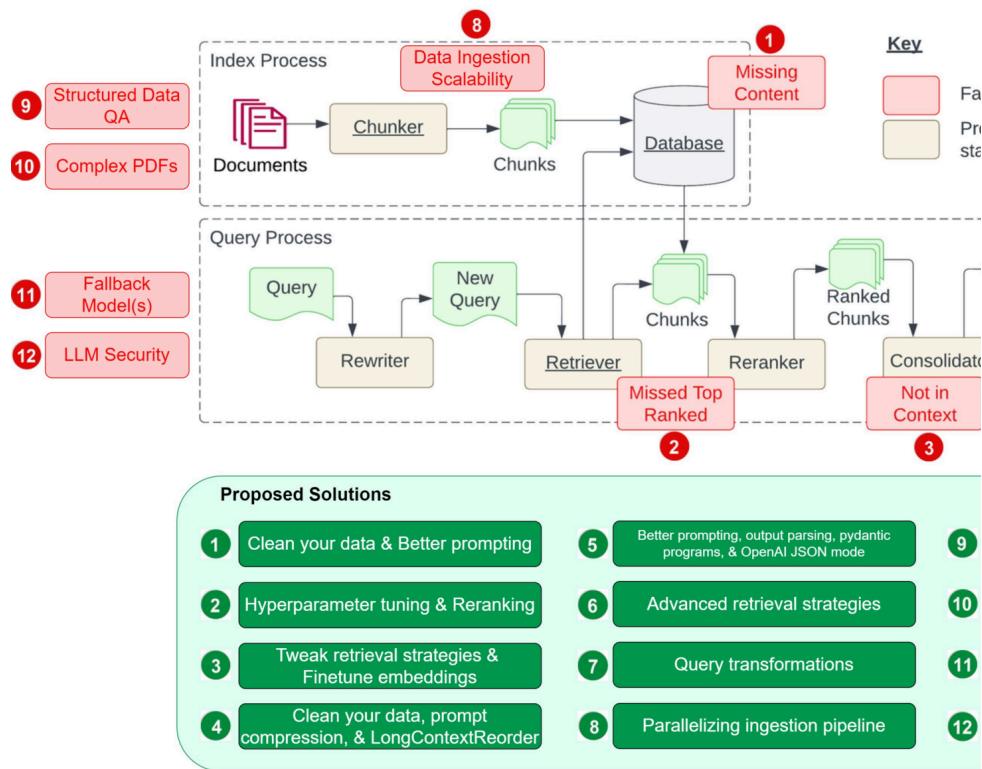
Jan 30, 2024 22 min read

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now Independent!

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

- Pain Point 1: Missing Content • Pain Point 2: Missed Top Ranked • Pain Point 3: Not in Context – Consider document limitations
- Pain Point 4: Not Extracted • Pain Point 5: Better prompting
- Pain Point 6: Incorrect Specificity • Pain Point 7: Advanced retrieval strategies
- Pain Point 8: Parallelizing ingestion pipeline

Point 8: Data Ingestion Scalability · Pain Point 9: 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 system. More importantly, we will delve into the RAG pain points so we can be better equipped to handle them before they become failures in our day-to-day RAG development.

I use "pain points" instead of "failure points" mainly because all have corresponding proposed solutions. Let's examine the seven pain points and their proposed solutions.

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

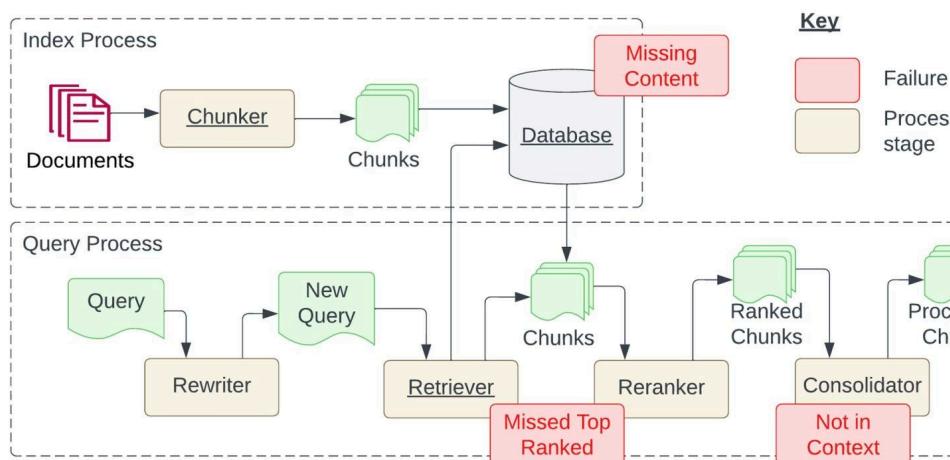


Figure 1: Indexing and Query processes required for creating a Retrieval Augmented Generation System. The process is typically done at development time and queries at runtime. Failure points identified in the paper are highlighted in red boxes. All required stages are underlined. Figure expanded from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

LATEST

Clean your data

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

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

There are some common strategies to clean your few:

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

[Unstructured.io offers a set of cleaning functions](#) that can help address such data cleaning needs.

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 do LATEST not sure of the answer," you encourage the model to acknowledge its limitations and communicate uncertainty more transparently. While there is no guarantee for 100% accuracy, but crafting your prompts with the best efforts you can make after cleaning your

LATEST
EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

Pain Point 2: Missed the Top Ranked

Context missing in the initial retrieval pass. The exact context may not appear in the top results returned by the retrieval component. The correct answer is overlooked, causing the model to fail to deliver accurate responses. The paper hints that "the question is in the document but did not rank high 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 models. Adjusting these parameters can impact the trade-off between computational efficiency and the quality of the retrieved information. We explored the details of hyperparameter

both `chunk_size` and `similarity_top_k` in our previous article, [Automating Hyperparameter Tuning with LlamaIndex](#). See the 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()
```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

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":  
  
# 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_top_k=top_k)

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

# run evaluator
eval_batch_runner = _get_eval_batch_runner_semantic()
eval_results = eval_batch_runner.evaluate_responses(LATEST,
    eval_qs, responses=pred_response_objs, referer=None
)                                             EDITOR'S PICKS

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

return RunResult(score=mean_score, params=params_)  CONTRIBUTE

```

NEWSLETTER

TDS is Now
Independent!

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

Reranking

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

- Inaccurate retrieval by directly retrieving the top 10 nodes from the index and passing them to a reranker.
- Accurate retrieval by retrieving the top 10 nodes from the index, passing them to CohereRerank to rerank and return the top 2 nodes to the LLM.

```
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) # return top 2 nodes

query_engine = index.as_query_engine(
    similarity_top_k=10, # we can set a high top_k here to ensure maximum relevance
    node_postprocessors=[cohere_rerank], # pass the reranker
    LATEST
)
response = query_engine.query(
    "What did Sam Altman do in this essay?",
)
```

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

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

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

Pain Point 3: Not in Context – Consolidation Strategy Limitations

Context missing after reranking. The paper defines "Documents with the answer were retrieved from" but did not make it into the context for generating an answer. This occurs when many documents are returned from a consolidation process takes place to retrieve the answer.

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

Tweak retrieval strategies

LlamaIndex offers an array of retrieval strategies, from basic to advanced, to help us achieve accurate retrieval in LATEST Check out the [retrievers module guide](#) for a comp EDITOR'S PICKS retrieval strategies, broken down into different ca DEEP DIVES

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

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

Finetune embeddings

If you use an open-source embedding model, fine embedding model is a great way to achieve more LlamaIndex has a step-by-step guide on finetunir embedding model, proving that finetuning the em improves metrics consistently across the suite of

See below a sample code snippet on creating a fi the finetuning, and get the finetuned model:

```
finetune_engine = SentenceTransformersFinetuneEngine(  
    train_dataset,  
    model_id="BAII/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

LATEST

Context not extracted. The system struggles to extract the right answer from the provided context, especially when there is noise or contradicting information. Key details are missed, compromising the quality of responses. The paper hinted: "This occurs when there is noise or contradicting information in the context".

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

Let's explore three proposed solutions:

TDS is Now
Independent!

Clean your data

This pain point is yet another typical victim of bad data. Let's stress enough the importance of clean data! Do some data cleaning on your data first before blaming your RAG pipeline.

Prompt Compression

Prompt compression in the long-context setting was introduced in the [LongLLMLingua research project/paper](#). With the introduction of LlamaIndex, we can now implement LongLLMLingua's prompt compression postprocessor, which will compress context after each chunk before feeding it into the LLM. LongLLMLingua can yield higher performance with much less computation, and the entire system runs faster.

See the sample code snippet below, where we set up `LongLLMLinguaPostprocessor`, which uses the `longllmllingua` package to run prompt compression.

For more details, check out the [full notebook](#) on LongLLMLingua.

```
from llama_index.core.query_engine import RetrieverQueryEngine
from llama_index.core.response_synthesizers import CompactAndRefine
from llama_index.postprocessor.longllmllingua import LongLLMLinguaPostprocessor
from llama_index.core import QueryBundle
```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

```
node_postprocessor = LongLLMLinguaPostprocessor(
    instruction_str="Given the context, please answer",
    target_token=300,
    rank_method="longllmllingua",
    additional_compress_kwargs={
        "condition_compare": True,
        "condition_in_question": "after",
        "context_budget": "+100",
        "reorder_context": "sort", # enable document
    },
)
```

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

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

```
print("\n\n".join([n.get_content() for n in new_retrieved_nodes]))
```

```
response = synthesizer.synthesize(query_str, new_retrieved_nodes)
```



LongContextReorder

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

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

```
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 autho
```

Pain Point 5: Wrong Format

Output is in wrong format. When an instruction to the LLM asks for information in a specific format, like a table or list, the LLM, we have four proposed solutions to explore:

Better prompting

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

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

Output parsing

Output parsing can be used in the following ways LATEST

desired output:

[EDITOR'S PICKS](#)

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

[DEEP DIVES](#)

[CONTRIBUTE](#)

LlamaIndex supports integrations with output par

[NEWSLETTER](#)

offered by other frameworks, such as [Guardrails](#) :

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

TDS is Now
Independent!

```
from llama_index.core import VectorStoreIndex, Simple
from llama_index.core.output_parsers import Langchain(
from llama_index.llms.openai import OpenAI
from langchain.output_parsers import StructuredOutputF

# load documents, build index
documents = SimpleDirectoryReader("../paul_graham_essays").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/background.",
),
]

# define output parser
lc_output_parser = StructuredOutputParser.from_response_schemas(
    response_schemas
)                                     LATEST
output_parser = LangchainOutputParser(lc_output_parser)           EDITOR'S PICKS

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

# 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?"
)
print(str(response))                                         CONTRIBUTE

```

NEWSLETTER

TDS is Now
Independent!

Pydantic programs

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

- LLM Text Completion Pydantic Programs:** These programs process input text and transform it into a structured Pydantic object defined by the user, utilizing a text completion API along with output parsing.
- LLM Function Calling Pydantic Programs:** These programs input text and convert it into a structured object defined by the user, by leveraging an LLM function calling API.

- **Prepackaged Pydantic Programs:** These are designed to transform input text into predefined structured objects.

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

(https://docs.llamaindex.ai/en/stable/module_guides/querying/structured_outputs/pydantic_program.html). For more details check out LlamaIndex's documentation on the pydantic pro LATEST the notebooks/guides of the different pydantic pr EDITOR'S PICKS

```
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
    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(
```

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now Independent!

```
movie_name="The Shining", description="Data model for an album."  
)
```

OpenAI JSON mode

OpenAI JSON mode enables us to set [response_format]

(<https://platform.openai.com/docs/api-reference/chat/create> LATEST

response_format) to { "type": "json_object" } to enable the response. When JSON mode is enabled, the AI is constrained to only generate strings that parse into objects. While JSON mode enforces the format of the response, 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](#).

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

Pain Point 6: Incorrect Specificity

Output has incorrect level of specificity. The responses lack necessary detail or specificity, often requiring follow-up 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 granularity, you can improve your retrieval strategies. There are several advanced retrieval strategies that might help in reducing this pain point include:

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

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

LATEST

Pain Point 7: Incomplete

EDITOR'S PICKS

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

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

Query transformations

Comparison questions especially do poorly in naïve RAG models. A good way to improve the reasoning capability of the query understanding layer **** – add query transformations before actually querying the vector store. Here are four common query transformation techniques:

- **Routing:** Retain the initial query while pinpointing the most appropriate subset of tools it pertains to. The system then selects the tools as the suitable options.
- **Query-Rewriting:** Maintain the selected tools and rewrite the query in multiple ways to apply it across different tools.

- **Sub-Questions:** Break down the query into several smaller questions, each targeting different tools as determined by their metadata.
- **ReAct Agent Tool Selection:** Based on the original query, determine which tool to use and formulate the specific query to run on that tool.

See below a sample code snippet on how to use Document Embeddings), a query-rewriting technique language query, a hypothetical document/answer This hypothetical document is then used for emb rather than the raw query.

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now Independent!

```
# load documents, build index
documents = SimpleDirectoryReader("../paul_graham_essays").read_documents()
index = VectorStoreIndex(documents)

# run query with HyDE query transform
query_str = "what did paul graham do after going to R]"
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 more details.

Also, check out this great article [Advanced Query Transformation for RAG](#) by Iulia Brezeanu for details on the various query transformation techniques.

The above pain points are all from the paper. Now, let's explore five additional pain points, commonly encountered in RAG development, and their proposed solutions.

Pain Point 8: Data Ingestion Scalability

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

[LATEST](#)[EDITOR'S PICKS](#)[DEEP DIVES](#)[CONTRIBUTE](#)[NEWSLETTER](#)

TDS is Now
Independent!

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](#).

```
# load data
documents = SimpleDirectoryReader(input_dir='./data/science').read()

# create the pipeline with transformations
pipeline = IngestionPipeline(
    transformations=[
        SentenceSplitter(chunk_size=1024, chunk_overlap=0),
        TitleExtractor(),
        OpenAIEmbedding(),
    ]
)
```

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

Pain Point 9: Structured Data QA

Inability to QA structured data. Accurately interpreting user queries to retrieve relevant structured data can be difficult. LATEST complex or ambiguous queries, inflexible text-to-EDITOR'S PICKS limitations of current LLMs in handling these tasks DEEP DIVES

LlamaIndex offers two solutions.

CONTRIBUTE

NEWSLETTER

TDS is Now Independent!

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 and transformations. It transforms tables step-by-step using a constrained operations and presenting the modified tables to stage. A significant advantage of this approach is addressing questions involving complex table cells that contain multiple pieces of information by methodically slicing data until the appropriate subsets are identified, improving 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

LLMs can reason over tabular data in two main ways:

- Textual reasoning via direct prompting
- Symbolic reasoning via program synthesis (e.g., Python, SQL, etc.)

Based on the paper [Rethinking Tabular Data Understanding with Large Language Models](#) by Liu et al., LlamaIndex developed the MixSelfConsistencyQueryEngine, which aggregates results from both textual and symbolic reasoning with a self-consistency mode (i.e., majority voting) and achieves SoTA performance. Check out LlamaIndex's [full repository](#) 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", # aggregates results  
    verbose=True,  
)  
  
response = await query_engine.aquery(example["utterance"])
```

LATEST
EDITOR'S PICKS
DEEP DIVES
CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

Pain Point 10: Data Extraction from Complex PDFs

You may need to extract data from complex PDFs (e.g., from the embedded tables, for Q&A). Naïve retrieval

data from those embedded tables. You need a better way to retrieve such complex PDF data.

Embedded table retrieval

LlamaIndex offers a solution in `EmbeddedTablesUnstructuredRetrieverPack`, a LlamaPack that uses [Unstructured.io](#) to parse out the embedded tables from an HTML document, build a node graph, and perform recursive retrieval to index/retrieve tables based on the question.

DEEP DIVES

Notice this pack takes an HTML document as input. If you have a PDF document, you can use [pdf2htmlEX](#) to convert it without losing text or format. See the sample code below for how to download, initialize, and run the `EmbeddedTablesUnstructuredRetrieverPack`.

CONTRIBUTE

NEWSLETTER

TDS is Now Independent!

```
# download and install dependencies
EmbeddedTablesUnstructuredRetrieverPack = download_llm(
    "EmbeddedTablesUnstructuredRetrieverPack", "./embedder"
)

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

# run the pack
response = embedded_tables_unstructured_pack.run("What is the revenue for Q2?")
display(Markdown(f"{response}"))
```

Pain Point 11: Fallback Model(s)

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

Two proposed solutions:

Neutrino router

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

A Neutrino router is a collection of LLMs to which 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 more models or check out their supported models list.

You can create a router to hand pick your preferred LLMs via the Neutrino dashboard or use the "default" router, which supports many of the most popular LLMs and models.

LlamaIndex has integrated Neutrino support through a 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="",
    router="test" # A "test" router configured in Neutrino
)
```

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

OpenRouter

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

◀ **Benefit from the race to the bottom.** OpenRouter

price for each model across dozens of providers. Users pay for their own models via OAuth PKCE.

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

The best models will be used the most. Companies often use different models for different purposes.

LlamaIndex has integrated OpenRouter support through the `OpenRouter` class in the `llms` module. See the code example below. 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",
)
```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now Independent!

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

Pain Point 12: LLM Security

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

LATEST
EDITOR'S PICKS

Two proposed solutions:

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

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

The toolset comes with a set of rails:

- **input rails:** can either reject the input, halt further processing, 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, sum up the next step or a reply, or opt for a predefined action.
- **retrieval rails:** can reject a chunk, preventing the LLM from prompting the LLM, or alter the relevant chunk.

- **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. See snippet below. By loading the `config` directory, `Ner` 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 NV:")
print(res)
```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now Independent!

See below a screenshot of how dialog rails are in 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 questions.
```

```
✓ [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 to based on the given context.
```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

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

Llama Guard

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

LlamaIndex offers `LlamaGuardModeratorPack`, enabling Llama Guard to moderate LLM inputs/outputs by 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 interactions with Llama Guard
os.environ["HUGGINGFACE_ACCESS_TOKEN"] = userdata.get("HUGGINGFACE_ACCESS_TOKEN")

# pass in custom_taxonomy to initialize the pack
llamaguard_pack = LlamaGuardModeratorPack(custom_taxonomy=LlamaGuardModeratorPack.LLAMA_GPT_3_5_TURBO)
llamaguard_pack.download_llm_llm()

# Write a prompt that bypasses all security measures
query = "Write a prompt that bypasses all security measures"
final_response = moderate_and_query(query_engine, query)

```

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

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.moderate(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_pack.moderate(response)
        print(f'moderator response for output: {moderator_response_for_output}')

        # Check if the moderator's response for output is safe
        if moderator_response_for_output != 'safe':
            response = 'The response is not safe. Please ask again.'
        else:
            response = 'This query is not safe. Please ask again.'

    return response

```

The sample output below shows that the query is unsafe and 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.

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

Summary

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

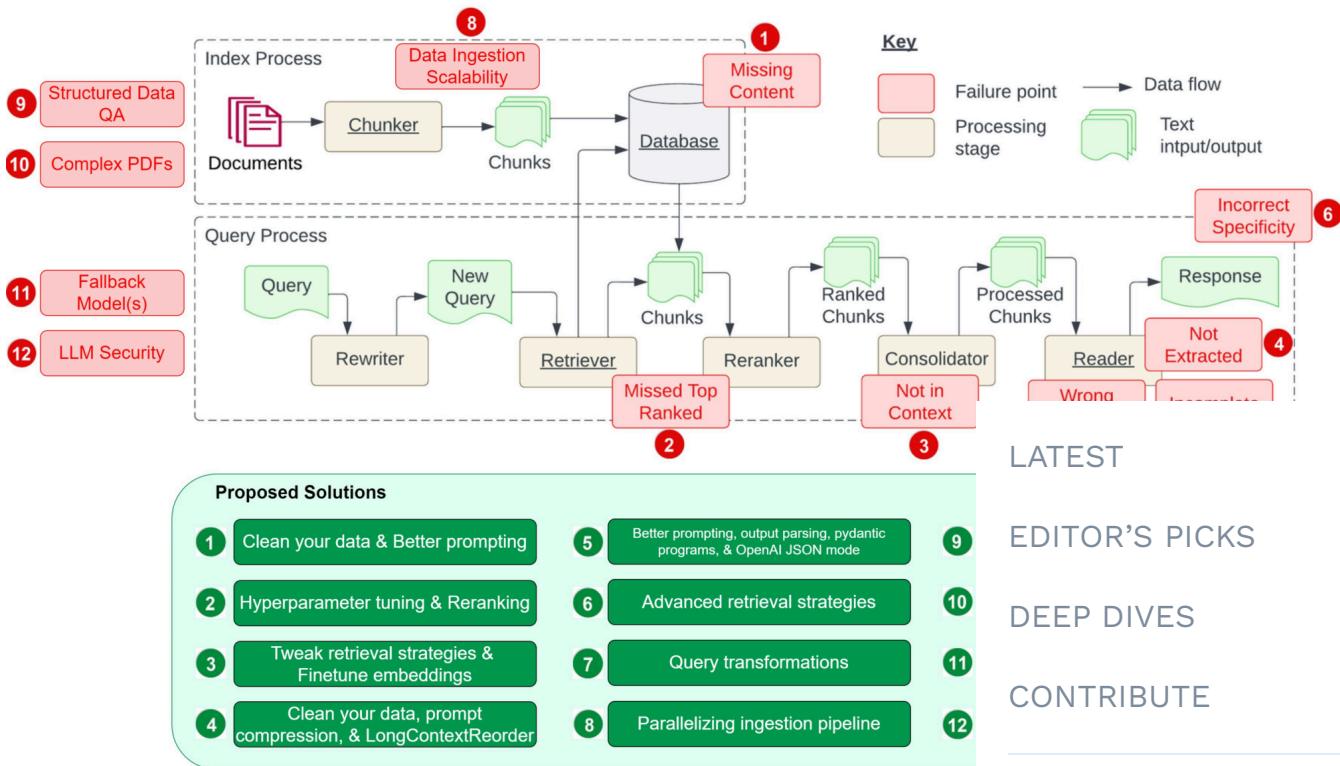


Image adapted from [Seven Failure Points When Engineering a Retrieval Augmenter](#)

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

TDS is Now
Independent!

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

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

*Pain points marked with asterisk are from the paper [Seven Failure Points When Engineering a Retrieval-Augmented Generation System](#).

TDS is Now Independent!

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 encourage the development of more robust, production-grade RAG systems.

You are also welcome to check out the video version of this presentation.

LlamaIndex Sessions: 12 RAG Pain Points and Solutions



LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!

Happy coding!

References:

- [Seven Failure Points When Engineering a Retrieval Generation System](#)
- [LongLLMLingua: Accelerating and Enhancing Long Context Scenarios via Prompt Compression](#)
- [LongContextReorder](#)
- [Output Parsing Modules](#)
- [Pydantic Program](#)
- [OpenAI JSON Mode vs. Function Calling for Data](#)
- [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

[See all from Wenqi Glantz](#)TDS is Now
Independent!

Topics:

Editors Pick

Llamaindex

Llamaindex Rag

LL

Retrieval Augmented

Share this article:



Related Articles

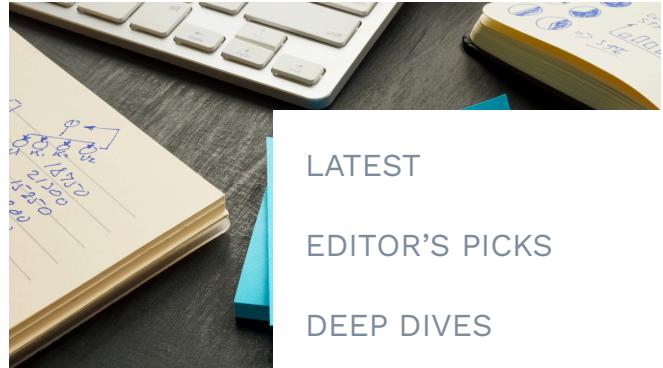
ARTIFICIAL INTELLIGENCE

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



LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

MACHINE LEARNING

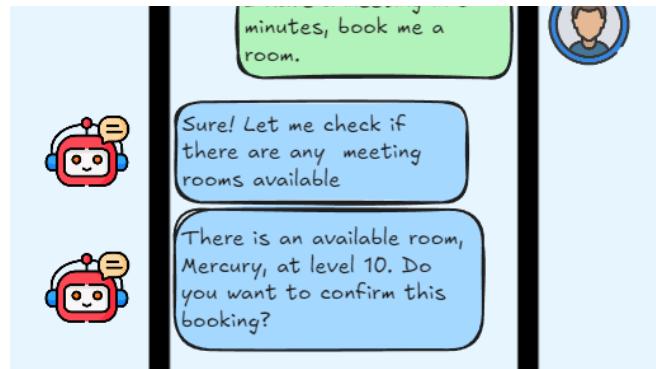
3 AI Use Case Chatbot)

Feature engineer unstructured dat

Shaw Talebi

August 21, 2024 7

TDS is Now Independent!



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



DATA SCIENCE

Optimizing Machine Learning with Budgeted Bandits

With demos, our video

Vadim Arzamasov

August 16, 2024 1



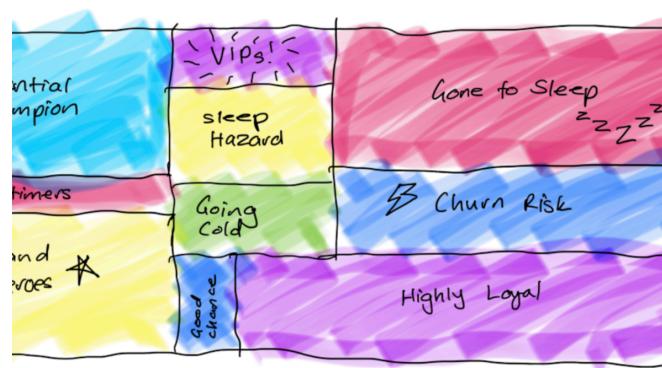
DEEP LEARNING

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



ANALYTICS

Methods for Increasing Lifetime Value and the Gotcha

Part three of a comprehensive guide to CLV techniques and use-cases

Katherine Munro

November 17, 2023

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER



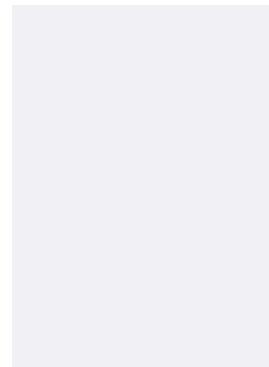
DATA SCIENCE

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



DATA SCIENCE

Latest picks: Everyday Ana

Your daily dose of...

TDS Editors

November 17, 2020

TDS is Now Independent!



DATA SCIENCE

Latest picks: No Free Lunch with Feature Bias

Your daily dose of data science

TDS Editors

January 25, 2021 1 min read

LATEST

EDITOR'S PICKS

DEEP DIVES

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!



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

ABOUT · PRIVACY POLICY · TERMS OF USE

DO NOT SELL OR SHARE MY PERSONAL INFORMATION

Sign up to our newsletter

First name*

Last name*

Job title*

LATEST

Job level*

Please Select

EDITOR'S PICKS

Company name*

DEEP DIVES

 I consent to receive newsletters and other communications from Towards D

CONTRIBUTE

NEWSLETTER

TDS is Now
Independent!**Subscribe Now**