# Data development for GenAI: A systems-level view
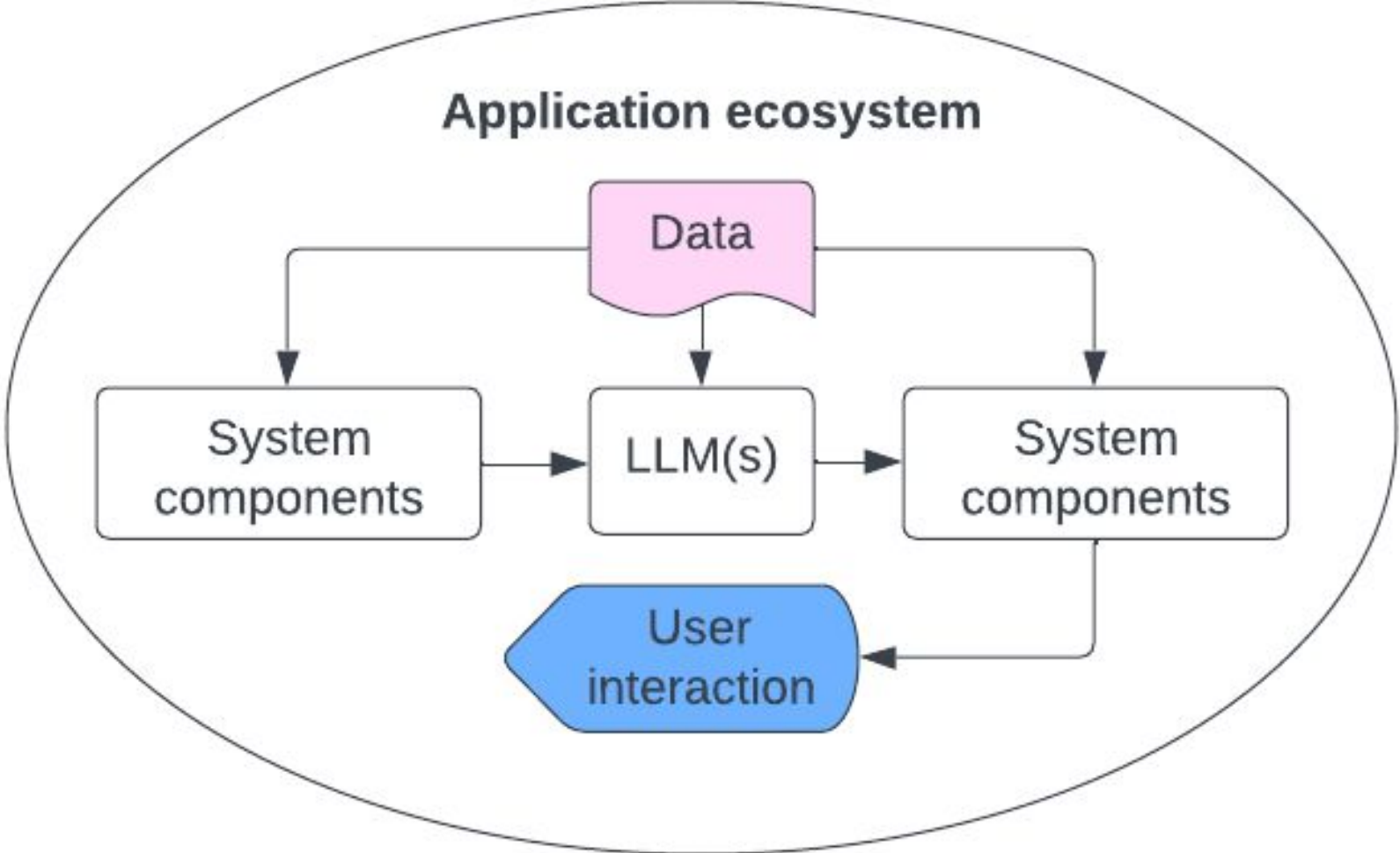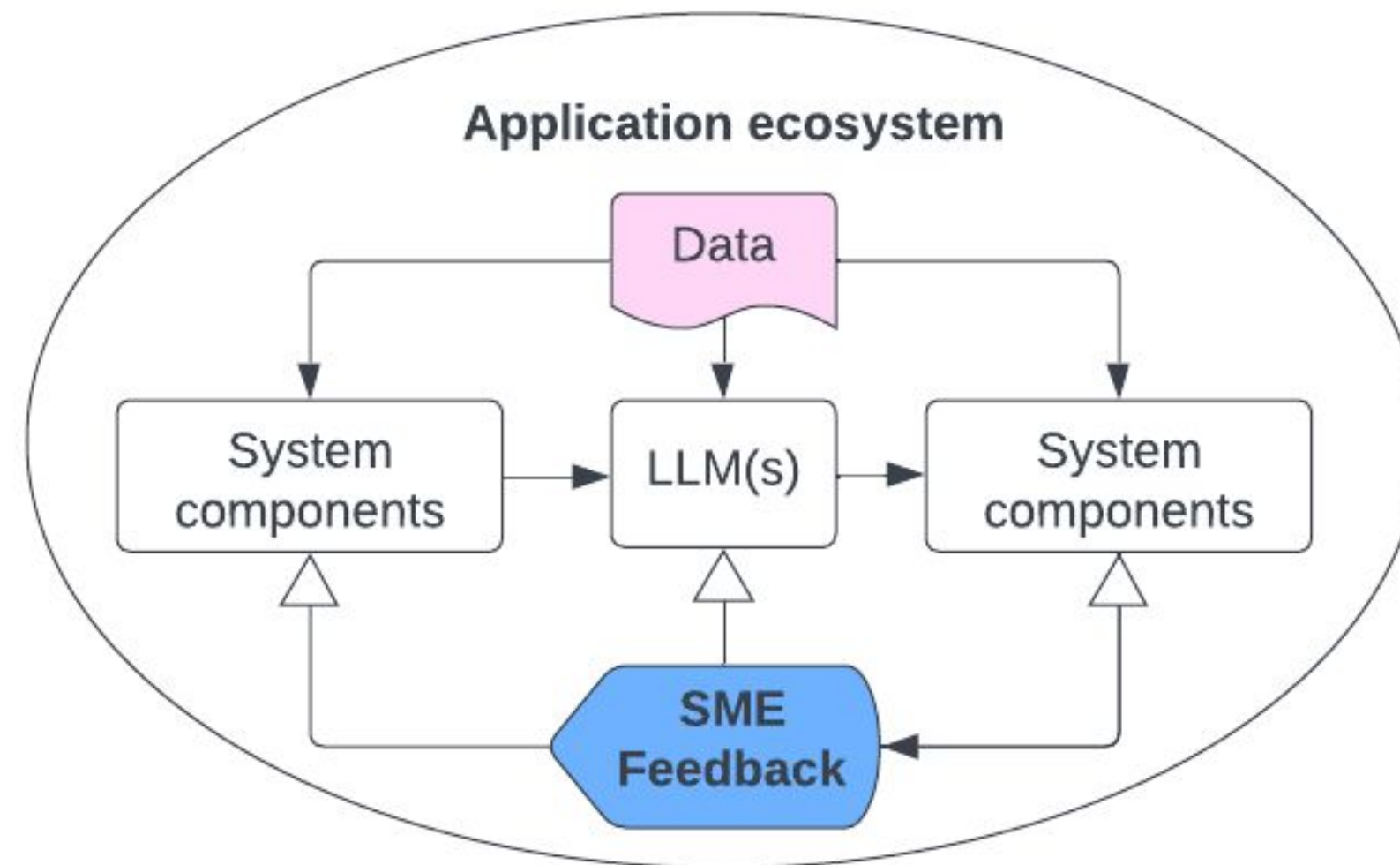
Chris Glaze
Staff Research Scientist

# Summary

➜ **For large language models (LLMs) in targeted business use cases:**
  - ◆ **Accuracy depends on the larger ecosystem in which they live**
  - ◆ **All system components benefit from fine-tuning with subject-matter expert (SME) feedback**

➜ **Snorkel is developing methods that efficiently incorporate SMEs in these development loops**
  - ◆ **Case Study: retrieval-augmented generation (RAG) for a top global bank, 54 point increase in question-answering accuracy in 3 weeks**

# Large language models do not exist in vacuums

# Many components require fine-tuning



*But subject-matter expert feedback has a scalability problem.*

# Keeping subject-matter experts in the loop

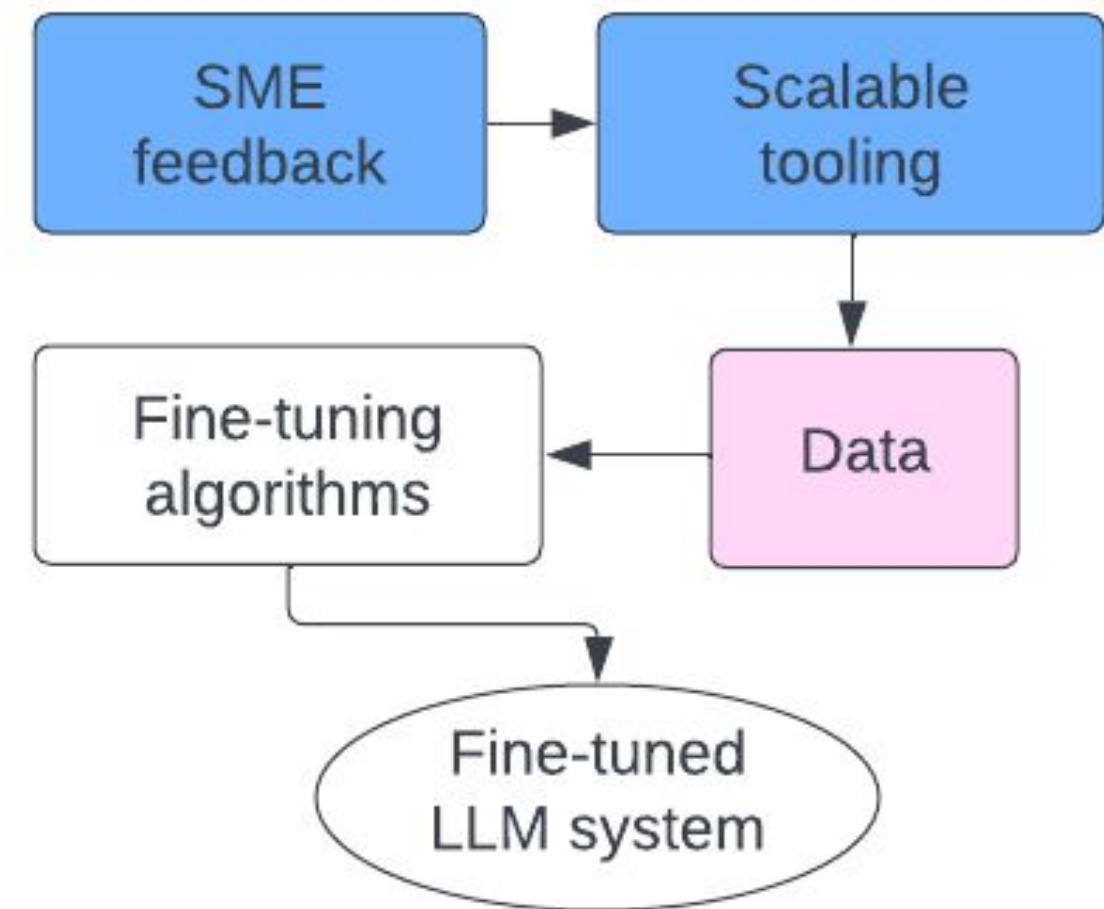**Snorkel thesis:**

*Data development is key.*

*Only subject-matter experts know what good looks like.*

**Snorkel approach:**

- **Keep subject-matter experts in the loop.**

- **Maximize value of their time with scalable methods to develop data.**
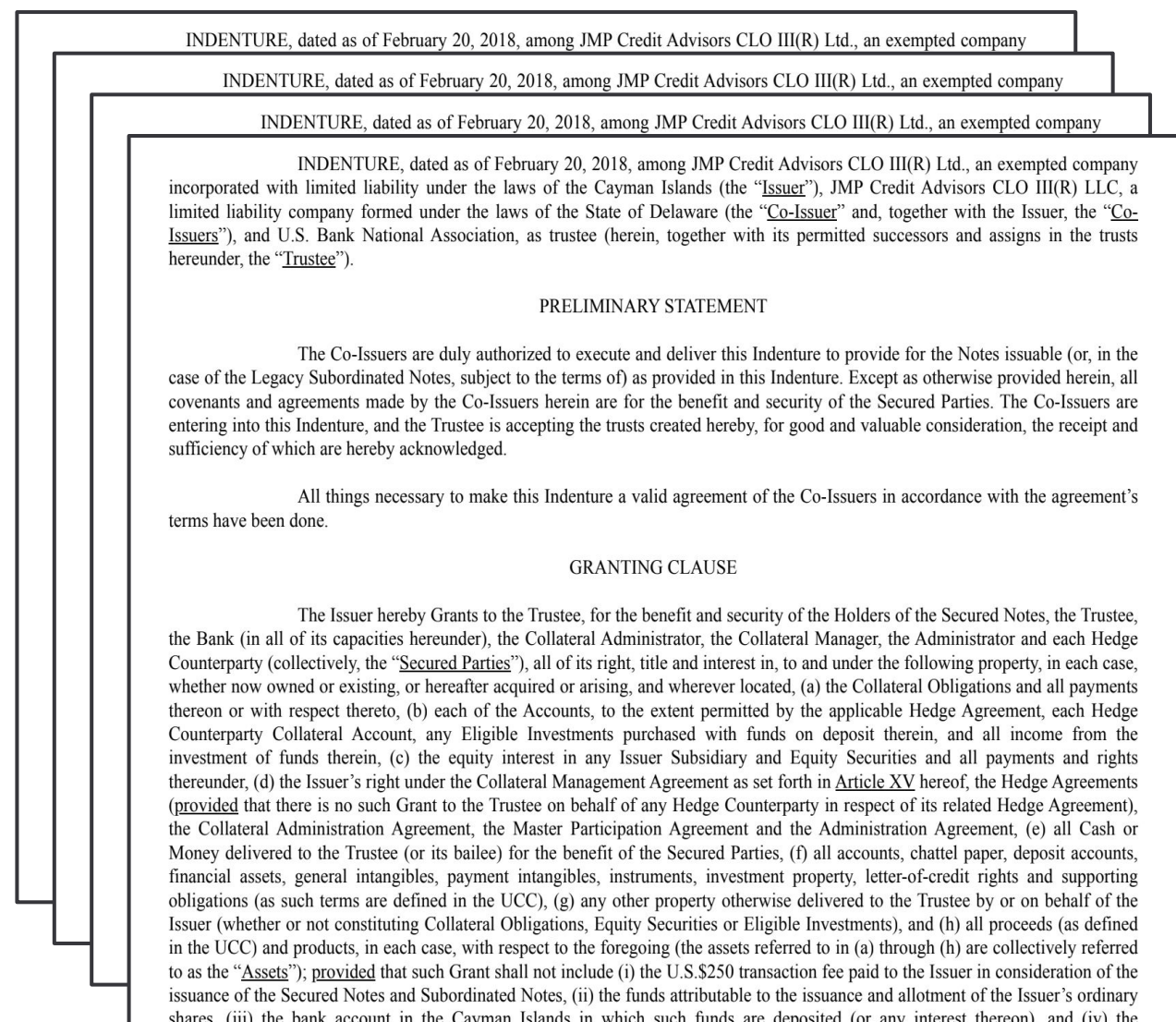
# Case Study: Retrieval-augmented generation (RAG) for a top global bank

We <u>fine-tuned</u> portions of an LLM-based question-answering system using <u>programmatic data development</u> techniques, over a <u>3 week period</u>:

| | Baseline LLM (GPT-4) + vector retrieval | **Fine-tuned** LLM Q&A system | *Improvement* |
|---|---|---|---|
| **Accuracy** | **25%** | **79%** | *+54 pts.* |

# Case Study: Retrieval-augmented generation (RAG) for a top global bank

## Unstructured Financial Documents
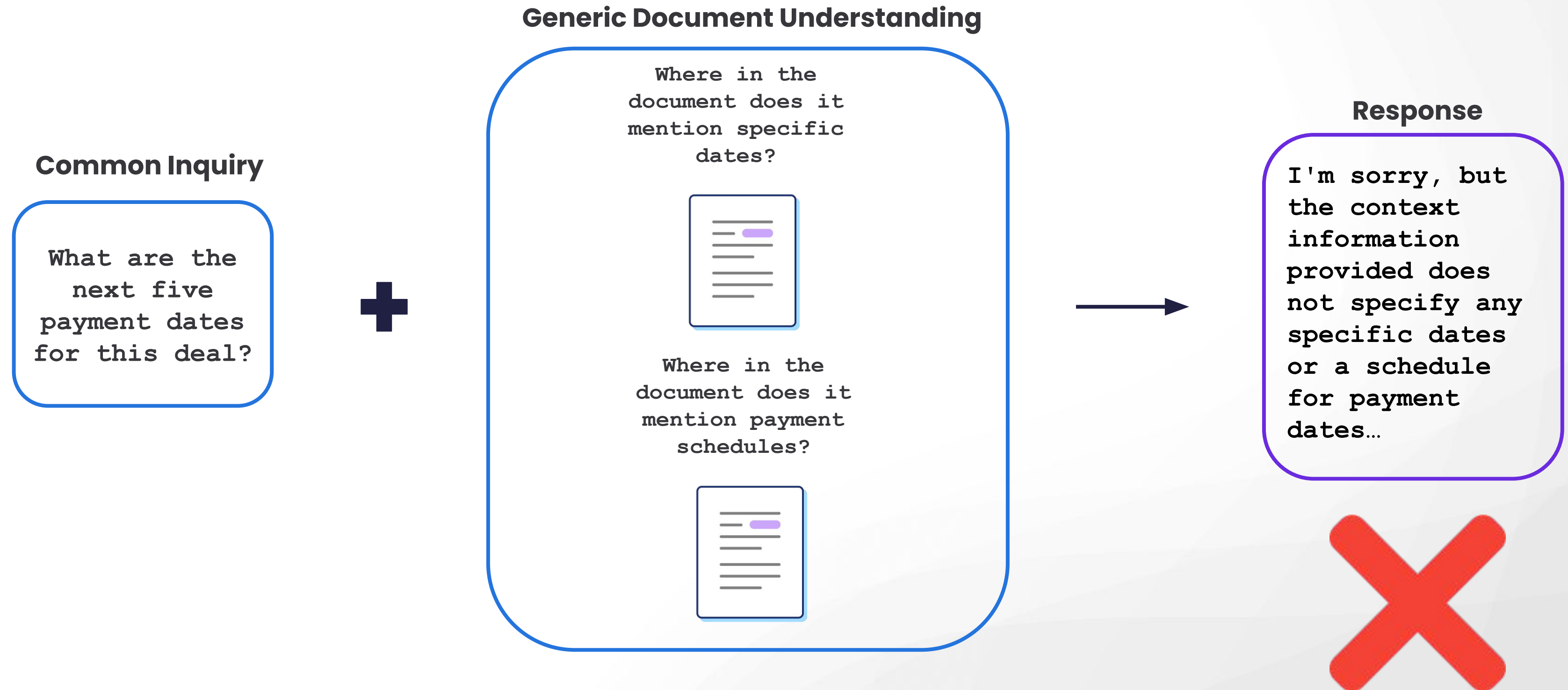


## Document characteristics

- **Domain Specific**: Content is specific toniche financial products

- **Dense Documents**: Each document ranges from 200-500+ pages long

- **Subject Matter Expertise**: Bank experts are required to comprehend content and produce accurate results

- **Complex business logic**: to extract relevant information from various parts of the document

# Example: Out-of-the-box RAG Performance

**Generic Document Understanding**

**Common Inquiry**

What are the next five payment dates for this deal?

**+**

Where in the document does it mention specific dates?

Where in the document does it mention payment schedules?

**→**

**Response**

I'm sorry, but the context information provided does not specify any specific dates or a schedule for payment dates…

# Example: Fine-tuned Performance

**Common Inquiry**

> What are the next five payment dates for this deal?

**+**

**Contextual Financial Document Understanding**

What is the jurisdiction and business days?

Page 18

What's the payment schedule?

Page 70

When is the first payment date?

Page 144

What is today's date?

External info

**+**

**Domain Knowledge**

Most recent payment + six business weeks.

Repeat five times OR until payment is complete

**Response**

> Payment dates for this deal are on a quarterly basis
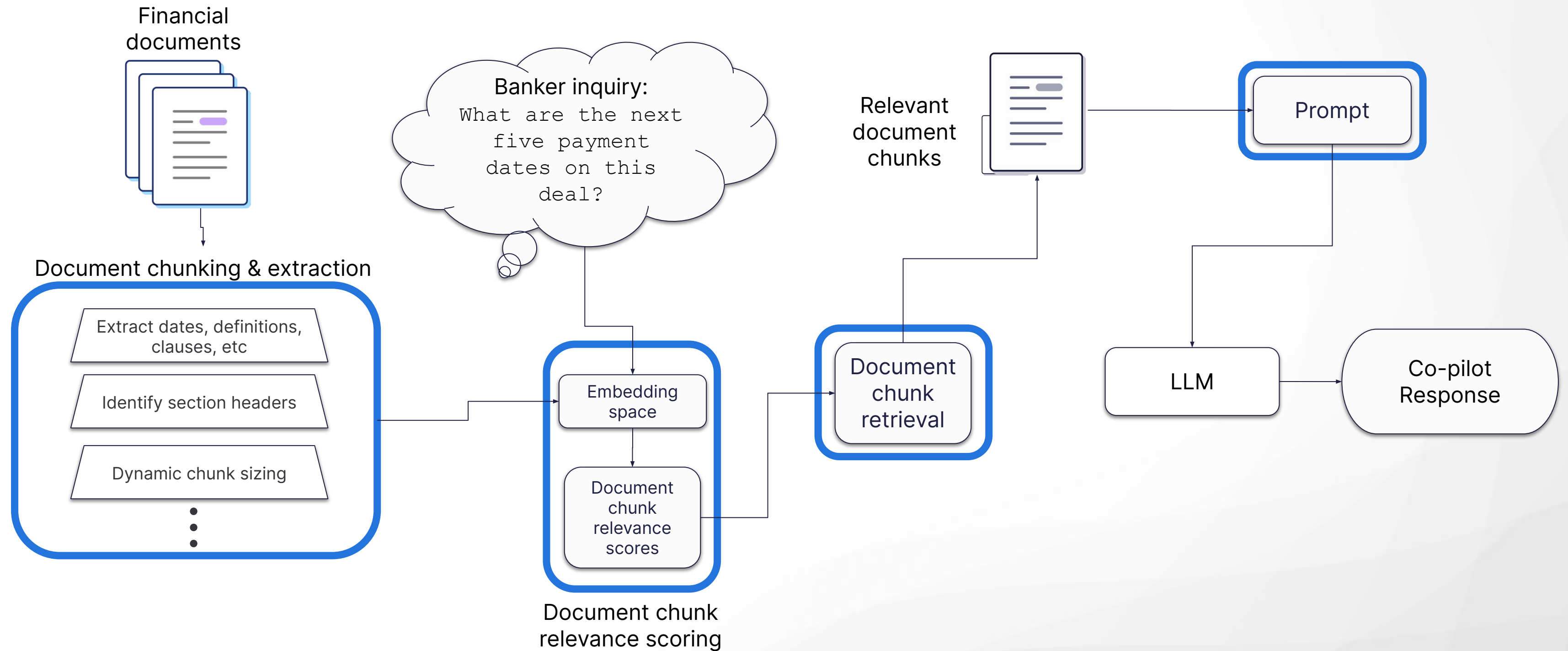> . . .
>
> The next five payment dates would be:
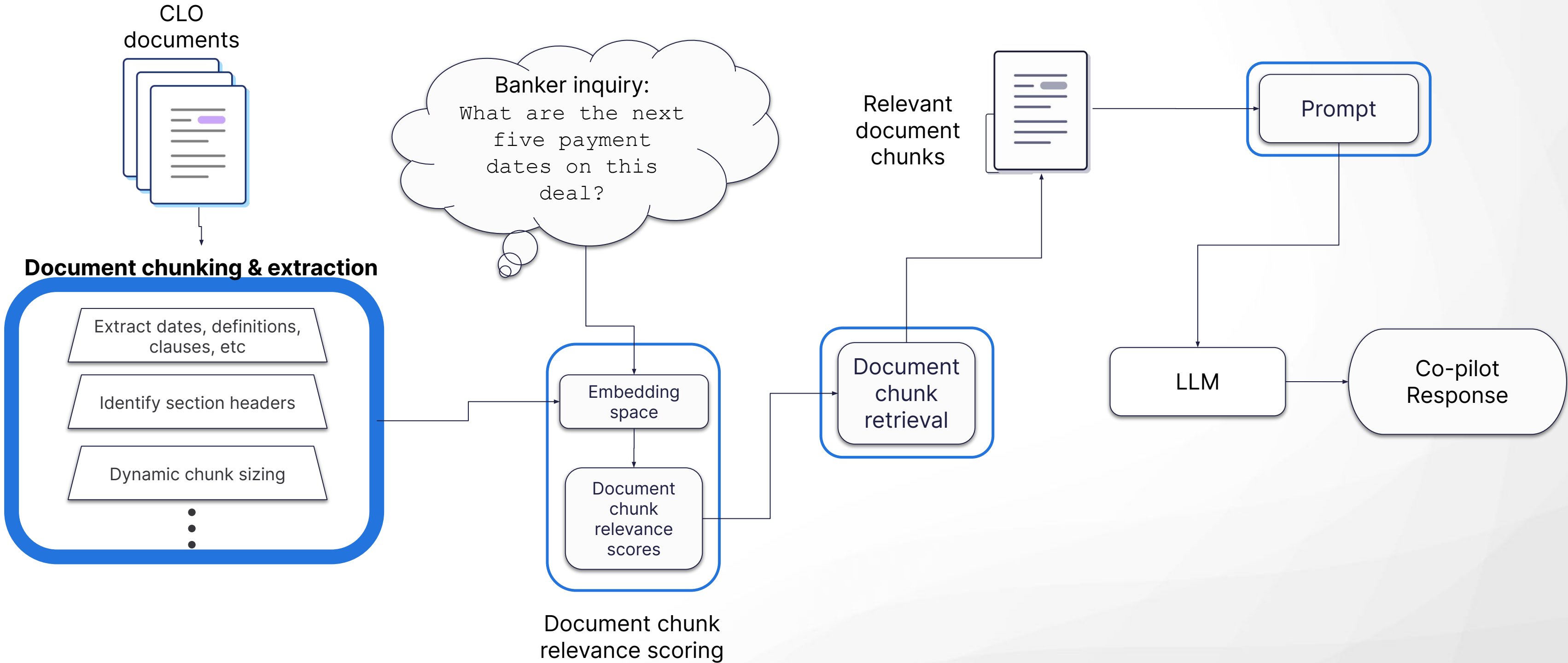> - **January 15, 2024**
> - **April 16, 2024**
>
> . . .

# How Snorkel Did It



*Fine tuned all components with subject-matter expert input → programmatic annotations*

# Document chunking, tagging, & extraction

CLO documents

**Document chunking & extraction**

Extract dates, definitions, clauses, etc

Identify section headers

Dynamic chunk sizing

Banker inquiry:
What are the next five payment dates on this deal?

Embedding space

Document chunk relevance scores

Document chunk relevance scoring

Document chunk retrieval

Relevant document chunks

Prompt

LLM

Co-pilot Response
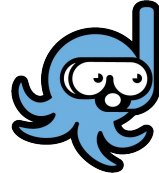
# Document tagging & extraction

**Challenge:** Q&A system struggled with key entities like dates, entities, definitions

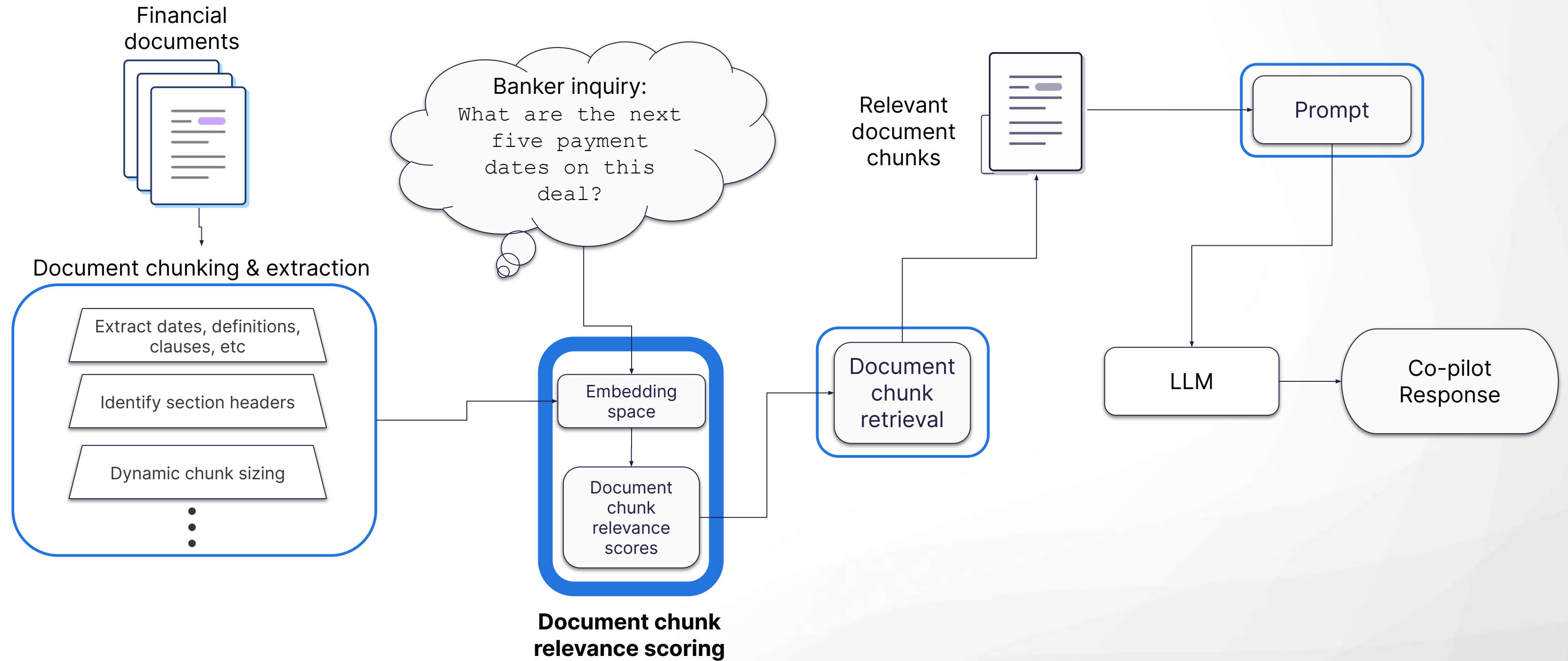**Solution:** Use Snorkel Flow to rapidly label structured entities

**Challenge:** Q&A system struggled to parse key sections, lumping them together

**Solution:** Use Snorkel Flow to parse section headers, tables, etc.

| | Time to build training data & ML Model | SF Generated Labels | Model Accuracy |
|---|---|---|---|
| Date Model | 4 hours | 141,000 spans | 99 F1 |
| Definitions Model | 4 hours | 43,100 spans | 93 F1 |

**Subject-matter expert manually annotated relevant sections and gave us logic for those annotations → retrieval of key information**

# Fine-tuning the relevance scoring model

Financial documents

Document chunking & extraction

Extract dates, definitions, clauses, etc

Identify section headers

Dynamic chunk sizing

Banker inquiry:
What are the next five payment dates on this deal?

Embedding space

Document chunk relevance scores

**Document chunk relevance scoring**

Document chunk retrieval

Relevant document chunks

Prompt

LLM

Co-pilot Response

# Fine-tuning the relevance scoring model
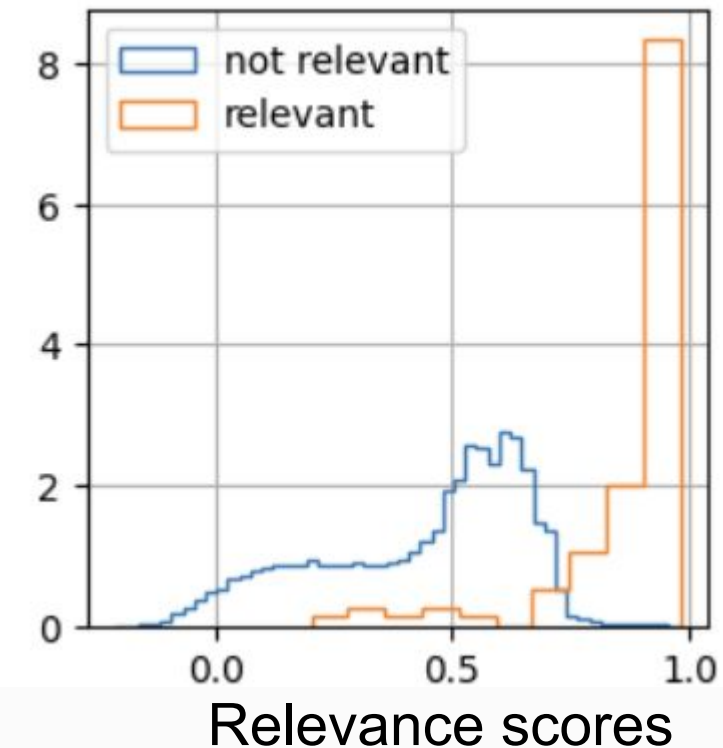


**Baseline pipeline**

- Pages and chunks ranked with Ada embedding
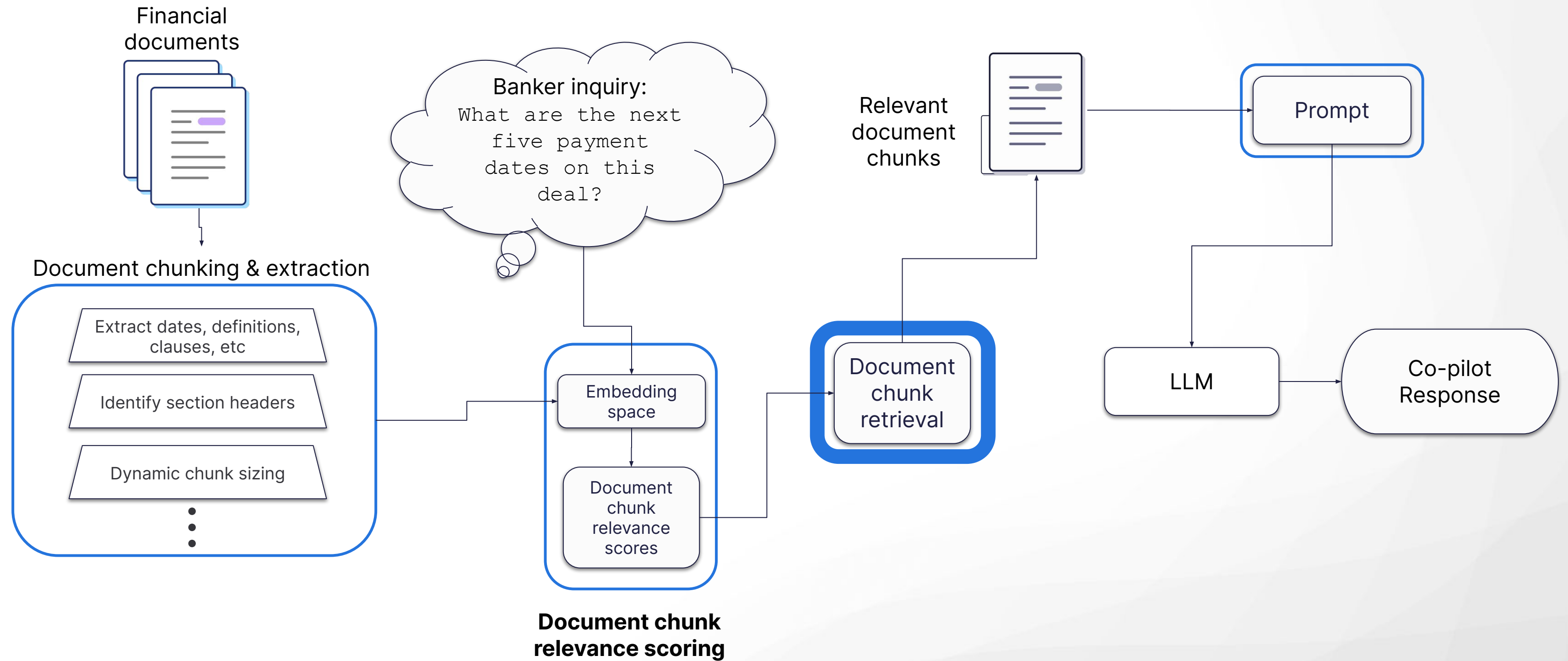
Relevance scores

**Fine tuned pipeline**

- Pages and chunks ranked with BGE embedding
- Embedding fine tuned with programmatic data

Relevance scores

**Challenge: Default embedding lumped together relevant and irrelevant chunks**

**Solution: Fine-tune embeddings to distinguished relevant and irrelevant chunks**
*Subject-matter expert annotations and logic used for training set*

# Fine-tuning the chunk retrieval algorithm

Financial documents

Document chunking & extraction

Extract dates, definitions, clauses, etc

Identify section headers

Dynamic chunk sizing

Banker inquiry:
What are the next five payment dates on this deal?

Embedding space

Document chunk relevance scores

**Document chunk relevance scoring**

Document chunk retrieval

Relevant document chunks

Prompt

LLM

Co-pilot Response
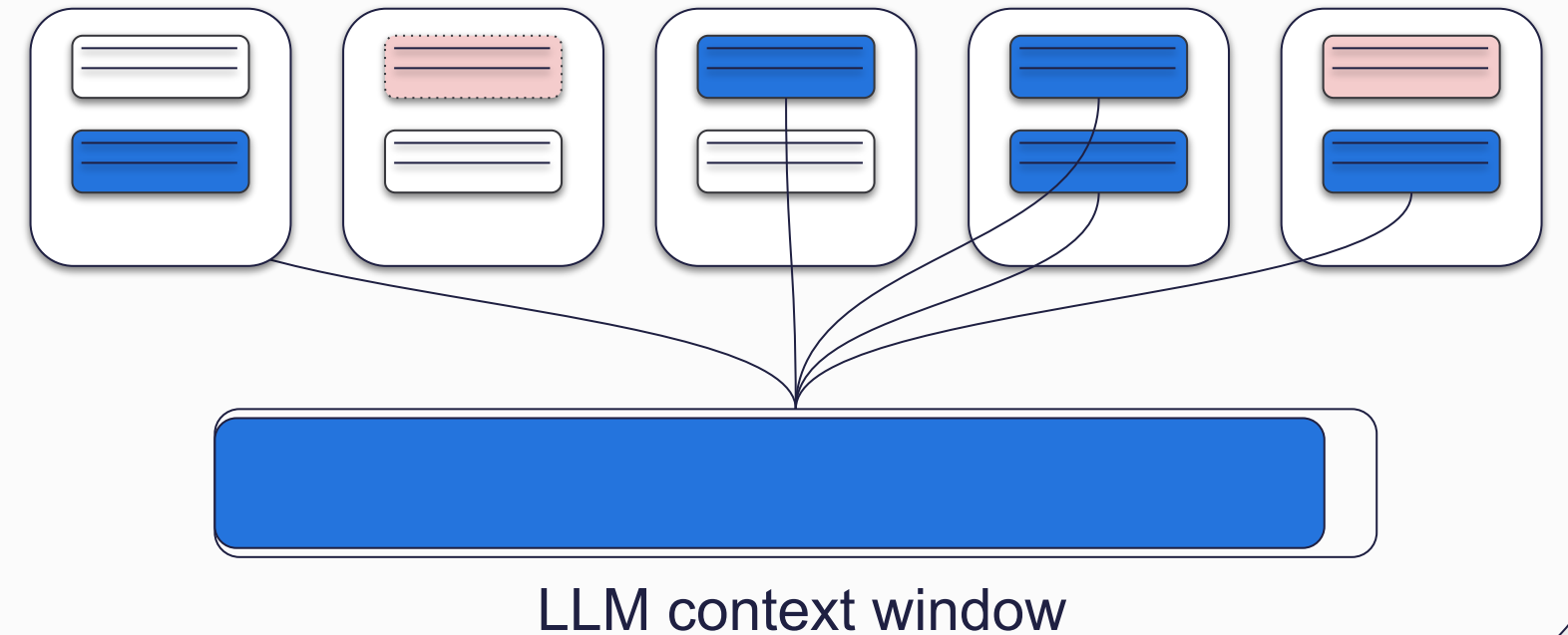
# Fine-tuning the chunk retrieval algorithm

## Baseline retrieval

- Fixed number (3) for every question based on ranking by relevance scores

## Fine-tuned retrieval

- Number of chunks adaptive, based on relevance scores and context window size

LLM context window
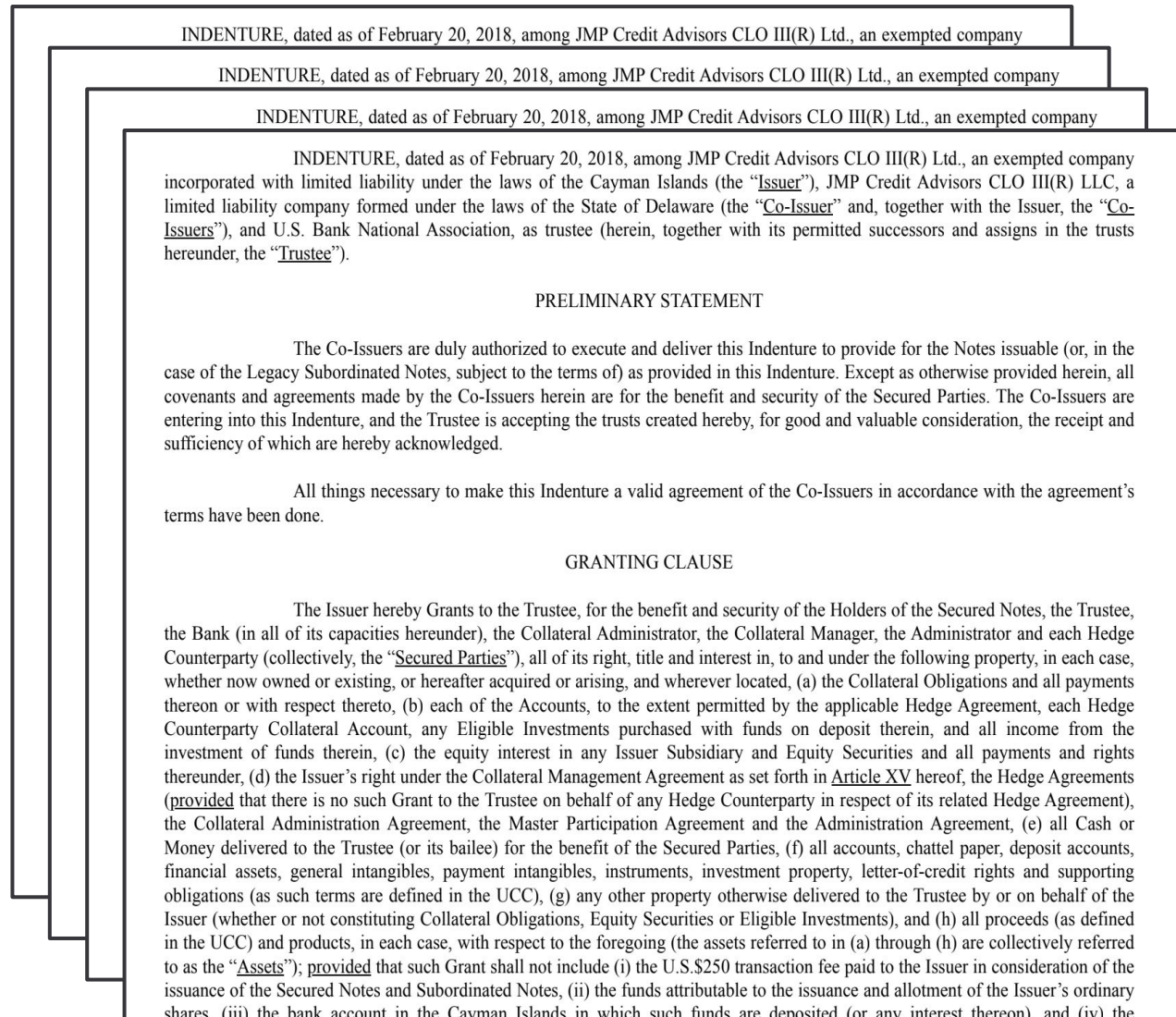
LLM context window

Low scoring chunks

High scoring chunks

**Challenge: Multiple, disparate document sections can be relevant to a given question**

**Solution: Allow for varying number of retrieved chunks based on relevance scores**
*Subject-matter expert annotations and logic used for training set*

# Results: 54 point increase in answer accuracy

**Unstructured Financial Documents**



INDENTURE, dated as of February 20, 2018, among JMP Credit Advisors CLO III(R) Ltd., an exempted company incorporated with limited liability under the laws of the Cayman Islands (the "Issuer"), JMP Credit Advisors CLO III(R) LLC, a limited liability company formed under the laws of the State of Delaware (the "Co-Issuer" and, together with the Issuer, the "Co-Issuers"), and U.S. Bank National Association, as trustee (herein, together with its permitted successors and assigns in the trusts hereunder, the "Trustee").

**25% Accurate**

With out-of-the-box LLM + RAG

*Failed to answer hardest questions*

**in 3 weeks**

**79% Accurate**

With fine-tuning + other data development

**Subject-matter expert**  **<10 hours**

**Ablation studies: ½ of increase came from fine tuning the embeddings, ½ from fine tuning other components**

# Summary

➜ **For large language models (LLMs) in targeted business use cases:**
- ◆ **Accuracy depends on the larger ecosystem in which it lives**
- ◆ **All system components benefit from fine-tuning with SME feedback**

➜ **Snorkel is developing methods that efficiently incorporate SMEs in these development loops**
- ◆ **Case Study: RAG for a top global bank, 54 point increase in QA accuracy in 3 weeks**

# Thank you!

**Chris Glaze**
**christopher.glaze@snorkel.ai**