# Welcome

## Build a Document Processing Pipeline for RAG Systems

### Stefan Krawczyk

Co-creator of Hamilton & Burr
Co-founder & CEO DAGWorks Inc.

maven

Components  - Code - Caveats

# Agenda: Build a Document Processing Pipeline for RAG Systems

- Components
- Code
- Caveats

maven

# Agenda: Build a Document Processing Pipeline for RAG Systems

- Components
- Code
- Caveats

**Leave you with:**

1. High level mental model of the process

2. Some code to help you get started

3. A sense for where caveats might lie on your journey

▥ maven

# Agenda: Build a Document Processing Pipeline for RAG Systems

- Components
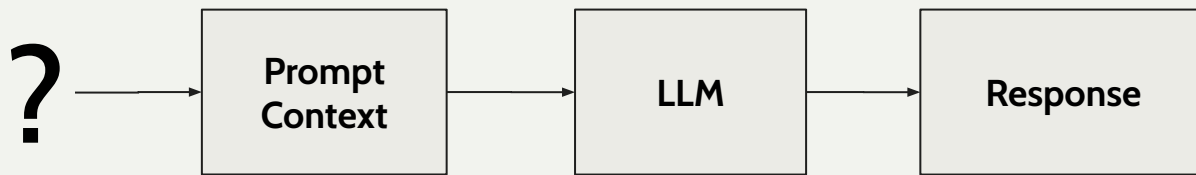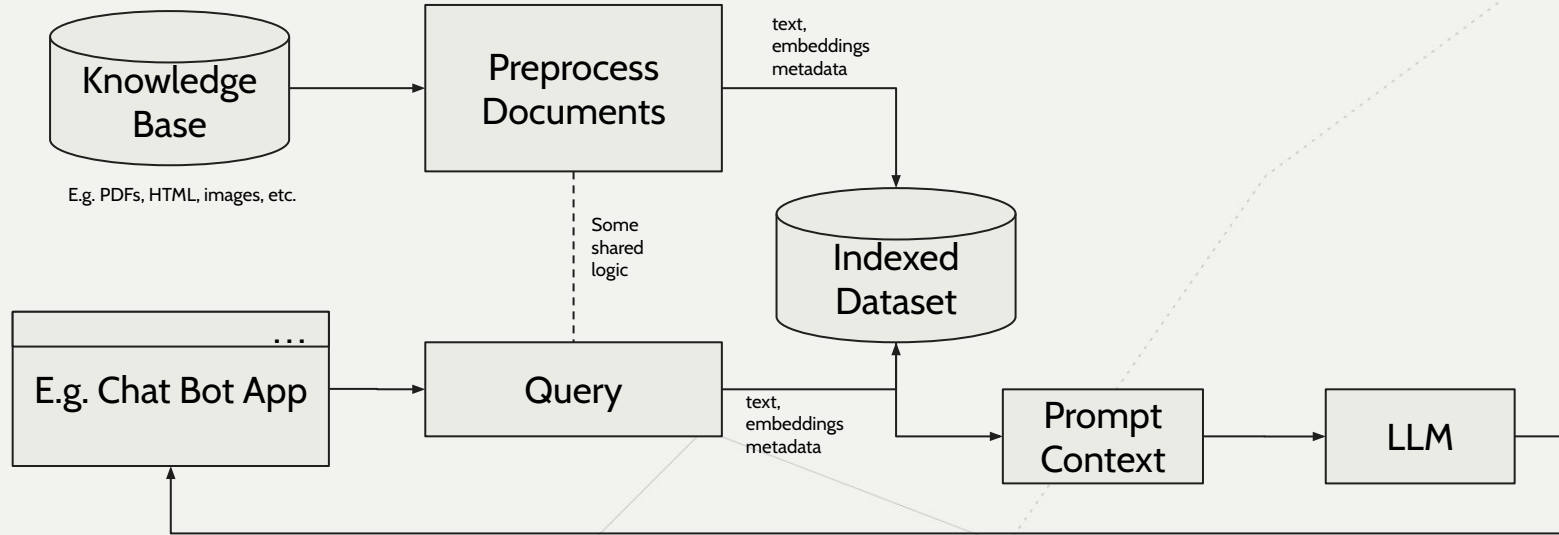- Code
- Caveats

In < 30 minutes!

**Leave you with:**

1. High level mental model of the process

2. Some code to help you get started

3. A sense for where caveats might lie on your journey
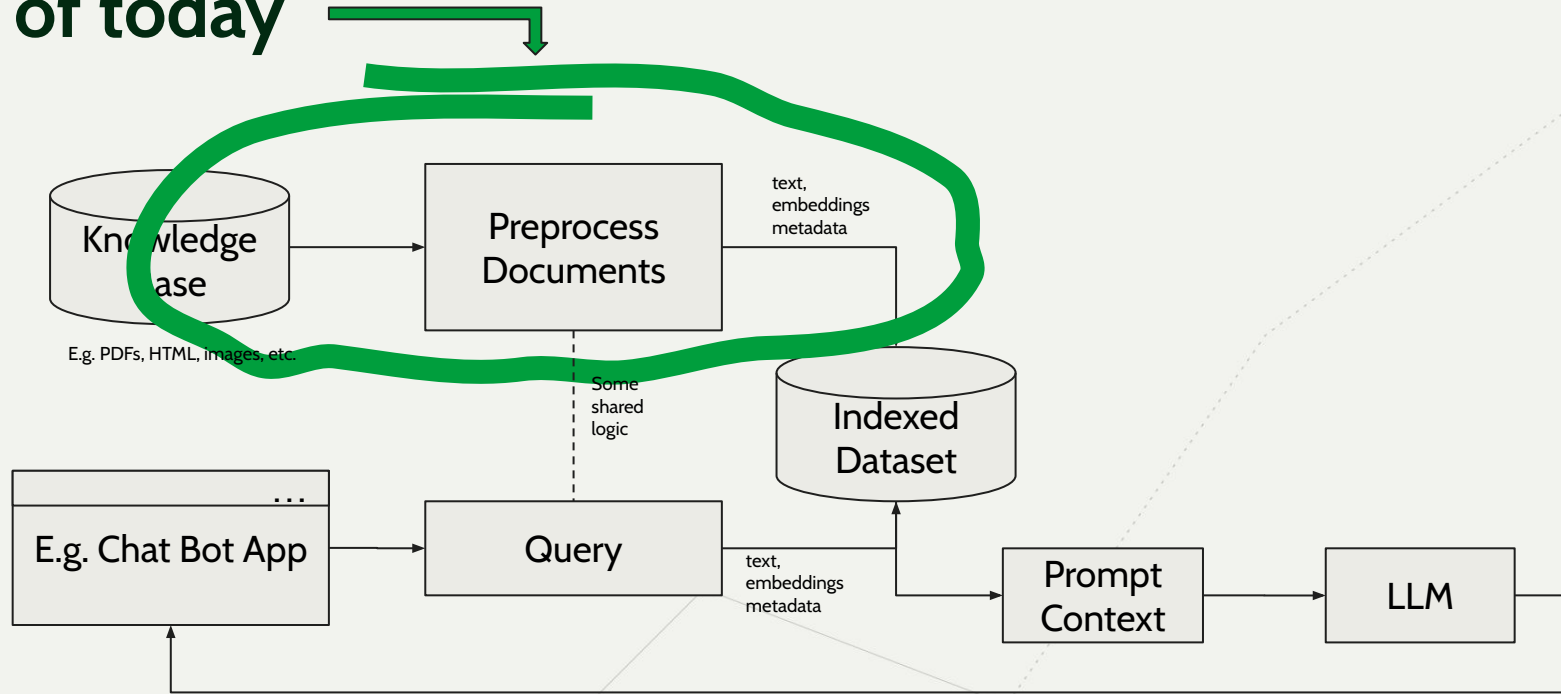
# Retrieval Augmented Generation



**?** → Prompt Context → LLM → Response

⇧

**Purpose:**
**Build the right context**
**for the LLM call**

maven

# Retrieval Augmented Generation



Knowledge Base

E.g. PDFs, HTML, images, etc.

Preprocess Documents

text, embeddings metadata

Some shared logic

Indexed Dataset

...

E.g. Chat Bot App

Query

text, embeddings metadata

Prompt Context

LLM

maven

# Focus of today



Knowledge base

E.g. PDFs, HTML, images, etc.

Preprocess Documents

text, embeddings metadata

Some shared logic

Indexed Dataset

E.g. Chat Bot App

Query

text, embeddings metadata

Prompt Context

LLM

maven

Big Picture

# Preprocessing Documents

| Load | → | Extract | → | Chunk | → | Embed | → | Store |
|------|---|---------|---|-------|---|-------|---|-------|

Get Document → Extract Text → Divide Text → Create Embedding → Store & index

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

Dolor sit amet consectetur adipiscing elit.

Id semper risus in hendrerit gravida rutrum quisque non tellus.

[0.2,0.3,0.4,0.5,0.2,]

[0.3,0.3,0.3,0.2,0.1,]

[0.5,0.4,0.2,0.8,0.7,]

**Note: 1 document → creates 1 or more embeddings**

maven

# Load

**Goal: get access to content for next steps**

Welcome to Hamilton

Hamilton is a general-purpose framework to write dataflows using regular Python functions. At the core, each function defines a transformation and its parameters indicates its dependencies. Hamilton automatically connects individual functions into a Directed Acyclic Graph (DAG) that can be executed, visualized, optimized, and reported on. Hamilton also comes with a UI to visualize, catalog, and monitor your dataflows.
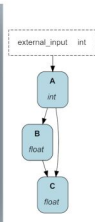
The ABC of Hamilton

## Why should you use Hamilton?

**Facilitate collaboration.** By focusing on functions, Hamilton avoids sprawling code hierarchy and generates flat dataflows. Well-scoped functions make it easier to add features, complete code reviews, debug pipeline failures, and hand-off projects. Visualizations can be generated directly from your code to better understand and document it. Integration with the Hamilton UI allows you to track lineage, catalog code & artifacts, and monitor your dataflows.

**Reduce development time.** Hamilton dataflows are reusable across projects and context (e.g., pipeline vs. web service). The benefits of developing robust and well-tested solutions are multiplied by reusability. Off-the-shelf dataflows are available on the Hamilton Hub.
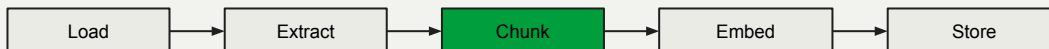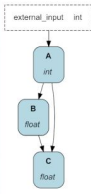
# Extract

**Goal: get text of interest**

**What you do here is** *very context* **dependent.**

**Spectrum:**

- **Simple:** rules, regular expressions, etc.

- **Middle:** OCR

- **Complex:** vision models





⊞ maven

Load → Extract → **Chunk** → Embed → Store

# Chunk

**Goal:**
- Segment for your use case

**Decisions:**
- How big are the chunks
- Is there overlap

**You'll use this:**
- to help find & build context for the LLM call.

Load → Extract → Chunk → **Embed** → Store

# Embed



The ABC of Hamilton

[0.384,0.417,0.122,0.176...

[0.680,0.305,0.598,0.415,0.824,0.204,0.115,0.504,0.145,0.086…]

[0.552,0.671,0.893,0.128,0.331,0.654,0.005,0.538,0.312,0.877,…]

**Goal:**
- Capture semantic meaning
- Used for finding similar text

**Decisions:**
- Implementation
  - E.g. openai, etc.
- Size
- Fine tuning or not.

# Store

Load → Extract → Chunk → Embed → **Store**

**Goal:** Store for retrieval

**What:** text, embeddings, metadata

**Decisions:** where do you store it?

[0.384,0.417,0.122,0.176,0.049,0.486,0.275,0.738,0.907,0.049...]

[0.680,0.305,0.598,0.415,0.824,0.204,0.115,0.504,0.145,0.086...]

[0.552,0.671,0.893,0.128,0.331,0.654,0.005,0.538,0.312,0.077,...]

Index & Store

# Example: Processing Hamilton's Documentation

Simple Pipeline Notebook (open in google collab)

# Caveats on the road to production

| Load | → | Extract | → | Chunk | → | Embed | → | Store |
|------|---|---------|---|-------|---|-------|---|-------|

→

→ Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

_

Dolor sit amet consectetur adipiscing elit.

_

Id semper risus in hendrerit gravida rutrum quisque non tellus.

→ [0.2,0.3,0.4,0.5,0.2,]

_

[0.3,0.3,0.3,0.2,0.1,]

_

[0.5,0.4,0.2,0.8,0.7,]

→

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

[0.2,0.3,0.4,0.5,0.2,]
_

Dolor sit amet consectetur adipiscing elit.

[0.3,0.3,0.3,0.2,0.1,]
_

Id semper risus in hendrerit gravida rutrum quisque non tellus

[0.5,0.4,0.2,0.8,0.7,]

**Get Document**

**Extract Text**

**Divide Text**

**Create Embedding**

**Store & index**

# Caveats on the road to production

Load → Extract → Chunk → Embed → Store

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

Dolor sit amet consectetur adipiscing elit.

Id semper risus in hendrerit gravida rutrum quisque non tellus.

[0.2,0.3,0.4,0.5,0.2,]

–

[0.3,0.3,0.3,0.2,0.1,]

–

[0.5,0.4,0.2,0.8,0.7,]

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

[0.2,0.3,0.4,0.5,0.2,]
–

Dolor sit amet consectetur adipiscing elit.

[0.3,0.3,0.3,0.2,0.1,]
–

Id semper risus in hendrerit gravida rutrum quisque non tellus

[0.5,0.4,0.2,0.8,0.7,]

**Get Document**  **Extract Text**  **Divide Text**  **Create Embedding**  **Store & index**

Domain specific

maven

# Caveats on

**Two main dimensions:**

- Domain specific
- Execution related

Load

Store

58  3  Norm and distance

common point. For example, the angle between the vectors $a = (1, 2, -1)$ and $b = (2, 0, -3)$ is

$$\arccos\left(\frac{5}{\sqrt{6}\,\sqrt{13}}\right) = \arccos(0.5661) = 0.9690 = 55.52°$$

(to 4 digits). But the definition of angle is more general; we can refer to the angle between two vectors with dimension 100.

The angle is a symmetric function of $a$ and $b$: We have $\angle(a,b) = \angle(b,a)$. The angle is not affected by scaling each of the vectors by a positive scalar: We have, for any vectors $a$ and $b$, and any positive numbers $\alpha$ and $\beta$,

$$\angle(\alpha a, \beta b) = \angle(a,b).$$

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

[0.2,0.3,0.4,0.5,0.2.]
–
Dolor sit amet consectetur adipiscing elit.

[0.3,0.3,0.3,0.2,0.1.]
–
Id semper risus in hendrerit gravida rutrum quisque non tellus

[0.5,0.4,0.2,0.8,0.7.]

**Get Document**

**Extract Text**

**Divide Text**

**Create Embedding**

**Store & index**

Domain specific

〽 maven

# Caveats on the road to production

E.g. rapid pace of change

E.g. Scale

Execution

| Load | → | Extract | → | Chunk | → | Embed | → | Store |

E.g. API vs GPU

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

Dolor sit amet consectetur adipiscing elit.

Id semper risus in hendrerit gravida rutrum quisque non tellus.

[0.2,0.3,0.4,0.5,0.2,]

–

[0.3,0.3,0.3,0.2,0.1,]

––

[0.5,0.4,0.2,0.8,0.7,]

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua.

[0.2,0.3,0.4,0.5,0.2,]

–

Dolor sit amet consectetur adipiscing elit.

[0.3,0.3,0.3,0.2,0.1,]

–

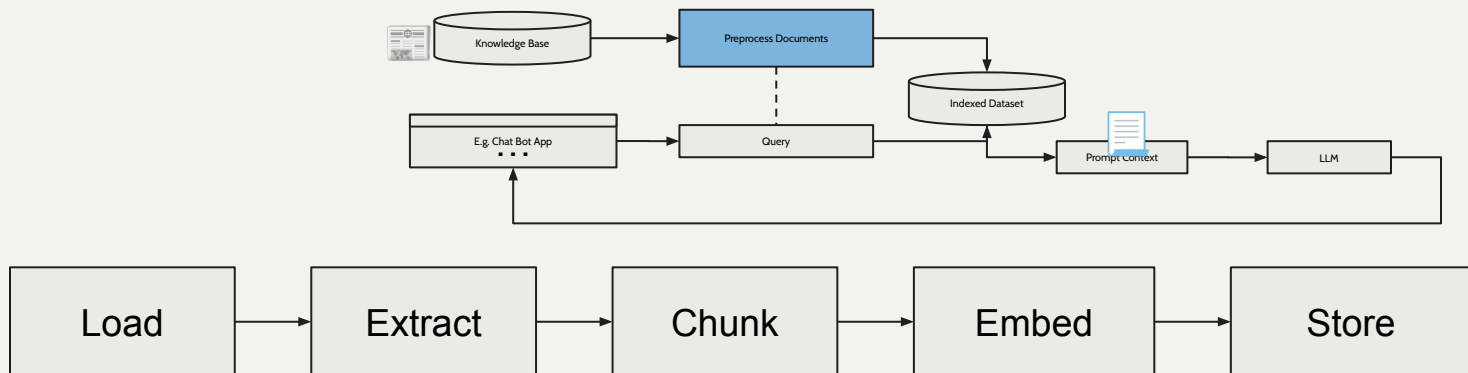Id semper risus in hendrerit gravida rutrum quisque non tellus

[0.5,0.4,0.2,0.8,0.7,]

**Get Document**

**Extract Text**

**Divide Text**

**Create Embedding**

**Store & index**

Domain specific

𝔐 maven

# Links / "Keywords" Slide

| | |
|---|---|
| Structuring Code | Hamilton (tryhamilton.dev) , Burr |
| Extracting Text | unstructured, OCR, LangChain, LlamaIndex, etc. |
| Chunking Text | unstructured, LangChain, LlamaIndex, etc. |
| Embedding Text | OpenAI, Anthropic, HuggingFace, etc. |
| Storage & Indexing | Files (e.g. parquet), numpy, PGVector, LanceDB, Marqo, etc |
| Scaling Processing | Ray, PySpark |
| More on embeddings | High-level article, More technical Google video |

▥ maven

# To finish



```
Load → Extract → Chunk → Embed → Store
```

**Thanks for listening!**

**Would you like a course on this?** Please fill out https://maven.com/forms/49388e

**Connect**: https://www.linkedin.com/in/skrawczyk/ (mention maven)

III maven