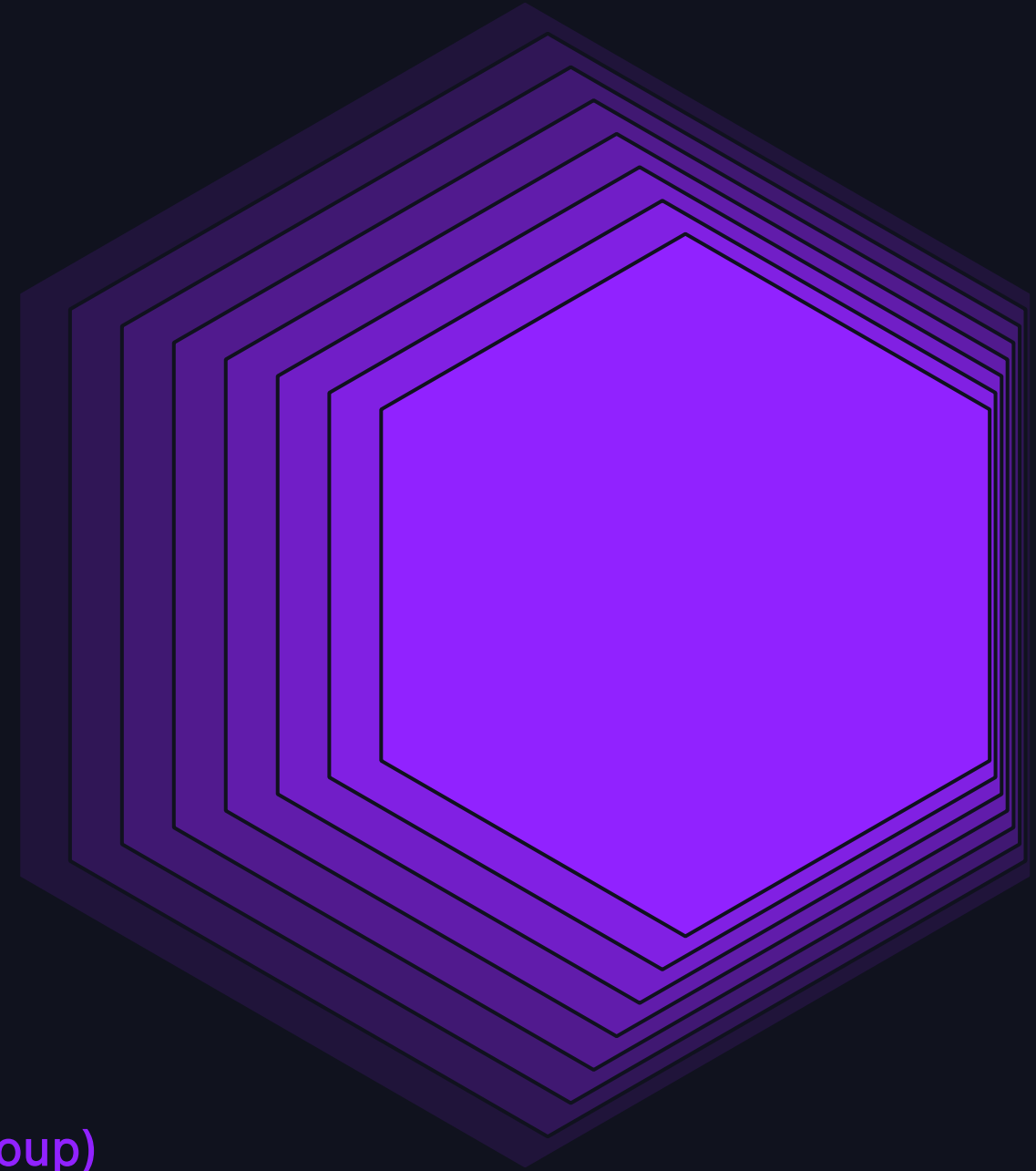


Augmenting Generative AI for complex Unstructured Data

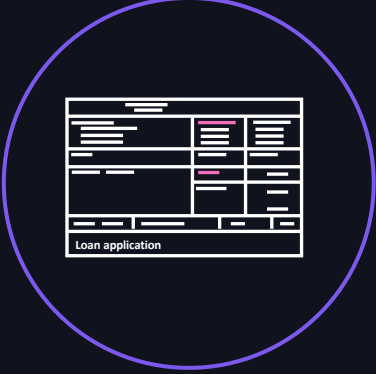


Venkat Viswanathan (AWS) and Mayur Rajdev (The Baldwin Group)
June 11th 2024

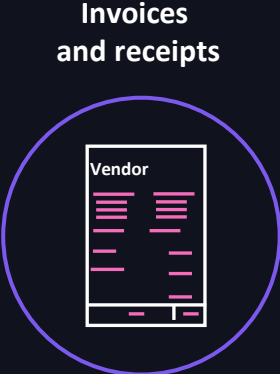
Complex Unstructured Data



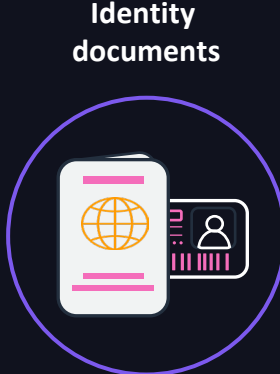
Complex PDFs



Forms



Invoices and receipts



Identity documents



Lending documents

Specialized documents



Handwriting



Tables



Signature

Industry Challenges when it comes to Complex Document Processing...

Manual Extraction

Time to Market

Classification

Scalability

Multi-Lingual Support

PII/PHI Security

Expensive

Decreased Employee Productivity

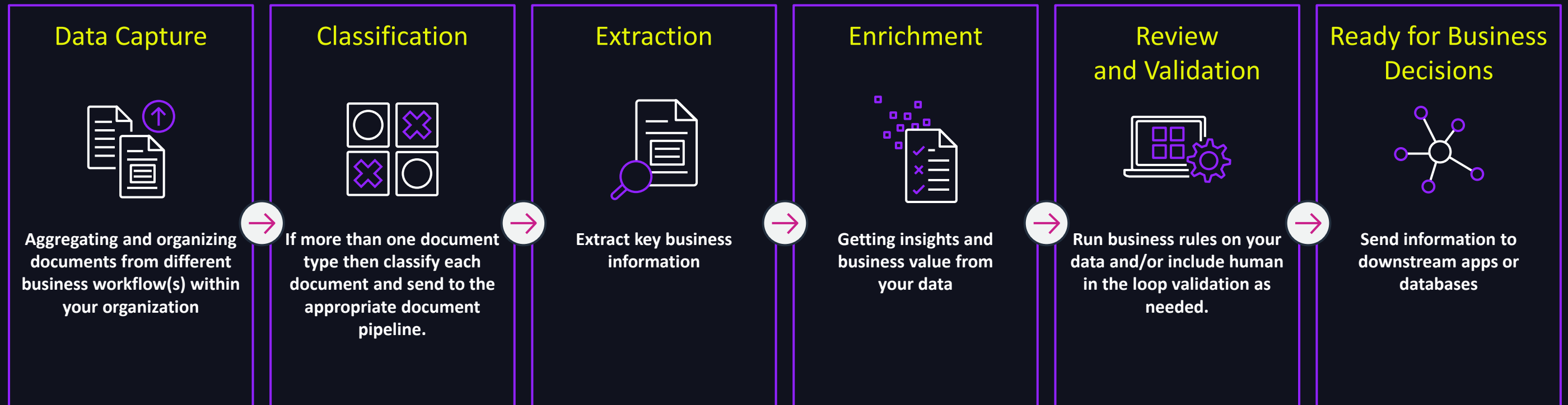
Inconsistent Output

Lack of Automation

Continuous Training

Language is not a barrier to run your business any more

Stages of Complex Document Processing



Business decision activities: Comparing documents with baseline, identifying gaps, and creating insights

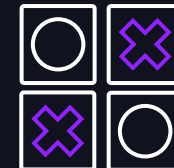


Key customer challenges



Data Capture

- Securely capture data
- Auto-correct for quality defects – distortion, dirt, rotated text, etc.



Classification

- Accurate identification and separation of documents where content across loan jackets may be highly variable



Extraction

- Capture data structures (tables, key-value pairs, entities, implied elements)
- Nested Data (page sections, duplicate fields)
- Data variations (SSN vs. Social Security #)



Enrichment

- Accurately identify PII (name, date of birth, addresses, etc.) out-of-the-box
- Flexibility to train custom models when necessary.



Validation

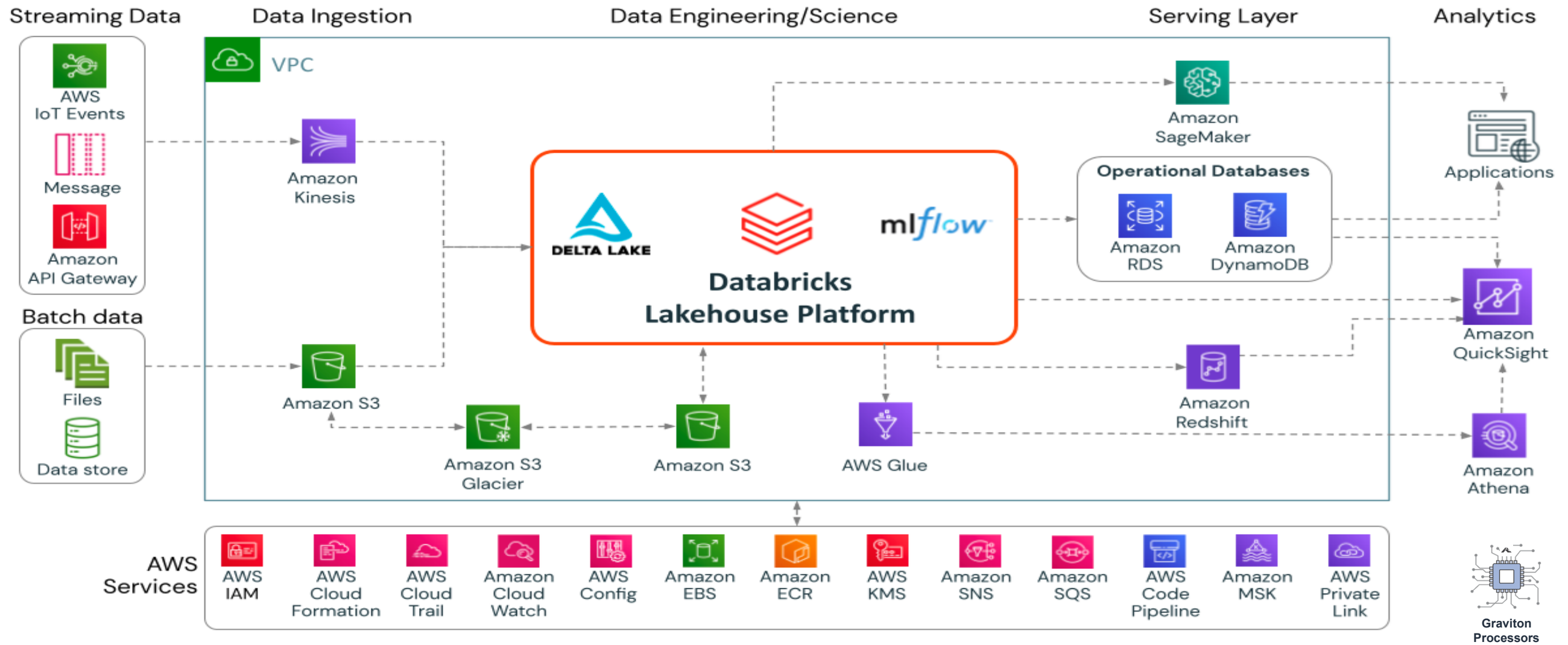
- Securely onboard your human review workforce to validate ML output
- Customize routing based on business rules, ML output metadata (e.g. confidence scores)



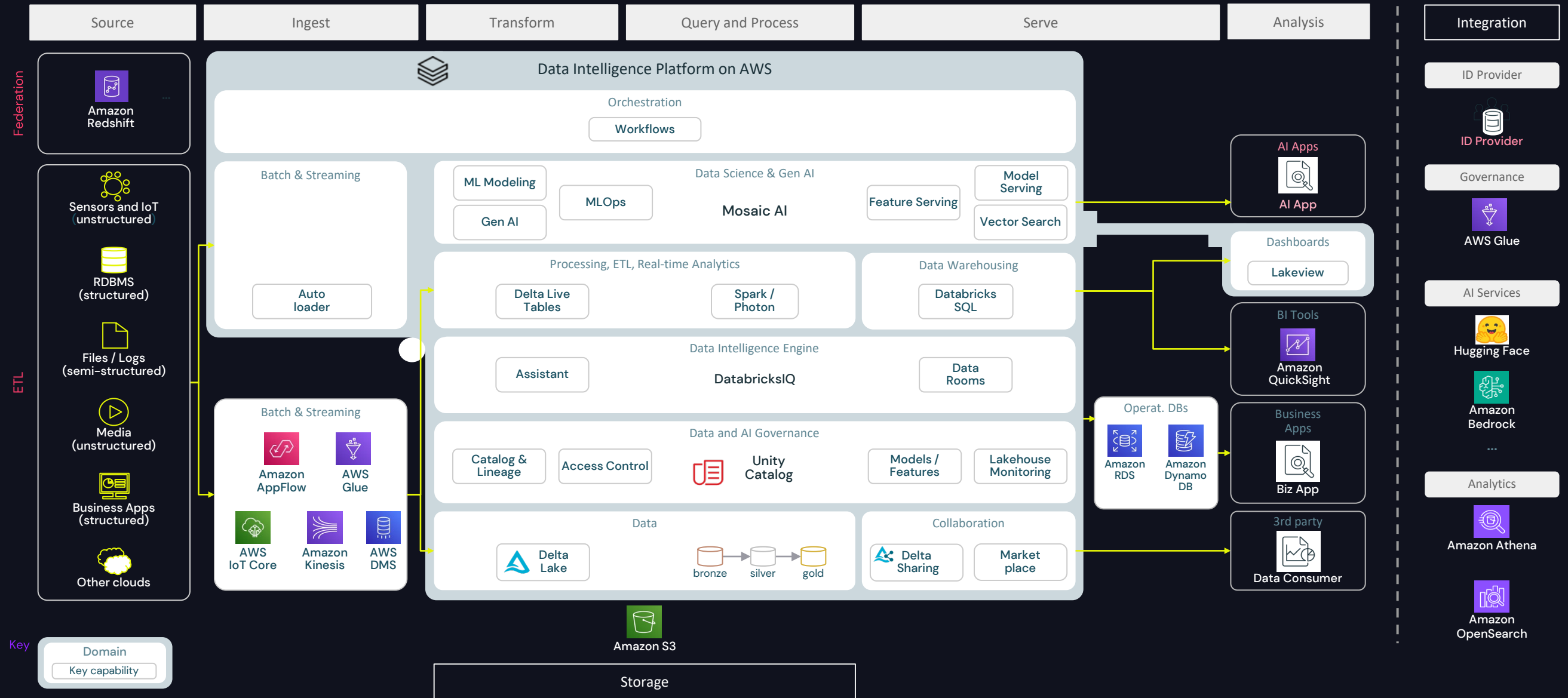
How AWS and Databricks together can help you speed up business processes, improve decision quality, and reduce overall costs with intelligent document processing (IDP)



Databricks on AWS: High-level Architecture



Databricks Data Intelligence Platform on AWS



Data Capture - Amazon S3

Aggregate and store documents securely



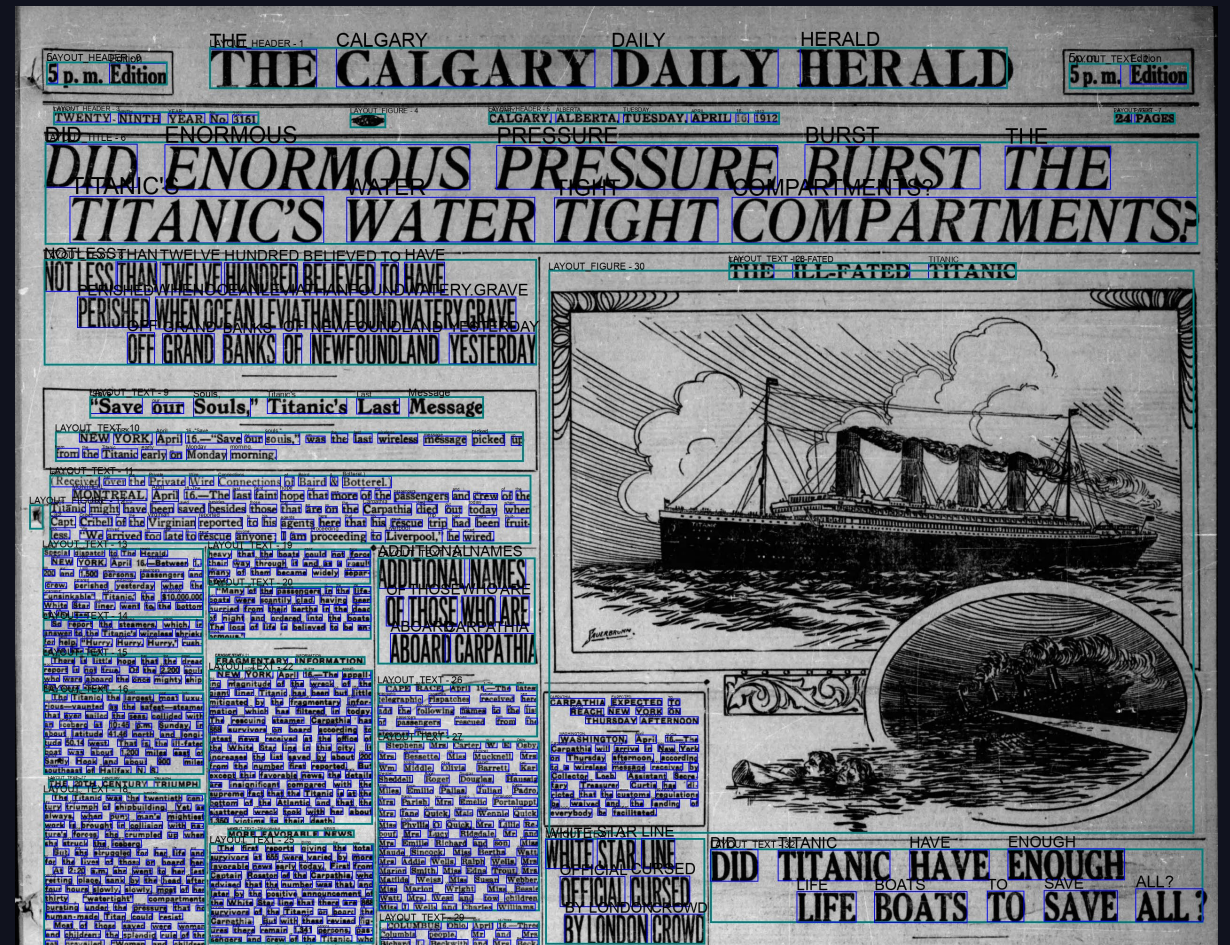
Automate document processing with Amazon Textract



Amazon Textract Response Objects

Layout Aware

- Main Title as **LAYOUT_TITLE**
- Text in the top margin of the document as **LAYOUT_HEADER**
- Text in the bottom margin of the document as **LAYOUT_FOOTER**
- Titles below the main title that represent sections in the document. Returned as **LAYOUT_SECTION_HEADER**
- Page number of the documents. as **LAYOUT_PAGE_NUMBER**
- Any information grouped together in list form as **LAYOUT_LIST**
- Location of an image in a document as **LAYOUT_FIGURE**
- Location of a table in the document as **LAYOUT_TABLE**
- Location of form key-value pairs in a document as **LAYOUT_KEY_VALUE**
- Text that is present typically as a part of paragraphs in documents as **LAYOUT_TEXT**



Extract data

Signature Detection

Detects Signatures in a document
(single page or multi page)

Captures page number where signature is
present

Also outputs confidence scores and
geometry info

Blocks:

PAGE, SIGNATURE

Example Output:

Signature Page 1 Accuracy 95%

Natural Language Query

Simple Natural Language Q & A Interface

Extract specific information of interest
from documents

Designed to work with just about any
type of document, no training needed

**Extract key entities from documents using natural
language questions**

Classify Document – Amazon Comprehend

Identify document types

Train a custom classifier to organize and identify different types of documents (e.g. W-2 paystubs, 1003, 1040, identity documents and more) within a loan jacket to inform downstream routing and application of custom business processing rules.

Training a custom classifier is as easy as providing 50 samples of each document type.

Real-time or asynchronous processing

Perform real time document classification when needed or reduce costs using asynchronous processes.



Classification



Entities
Default & custom



PII
Personally Identifiable
Information

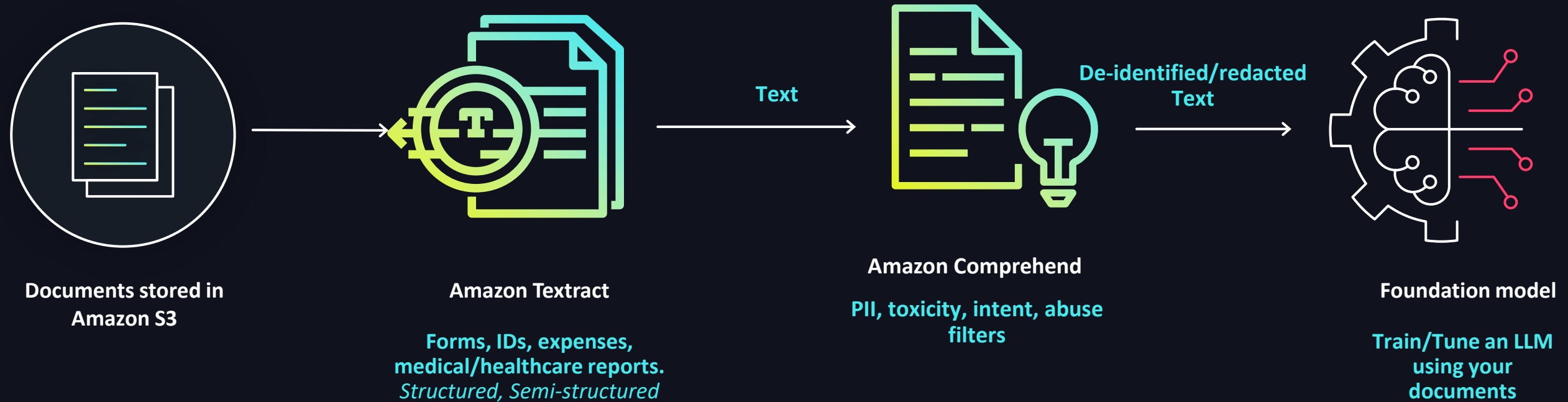


Key phrases



Sensitive data handling for generative AI use-cases

Redacting inputs can improve sensitive data handling with Generative AI models.



Redaction of sensitive data

DOCUMENT TEXT

Hi, my name is John Doe. For verification, the last 4 digits of my social are 6789 and my DOB is 01/01.



Hi, my name is [NAME]. For verification, the last 4 digits of my social are [SSN] and my DOB is [DATE_TIME].

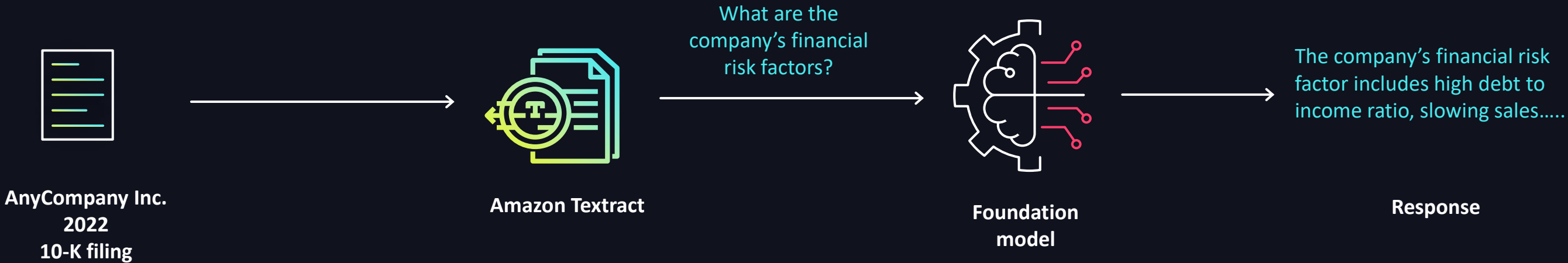


Hi, my name is *****. For verification, the last 4 digits of my social are **** and my dob is *****.



Generative AI - Zero-shot and few-shot prompts

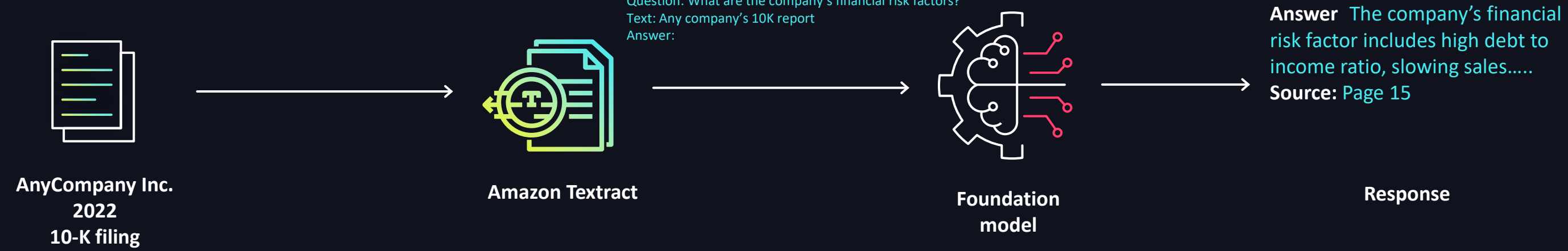
Get the desired output from an LLM with prompt engineering



Answer the question given the text and include sources.

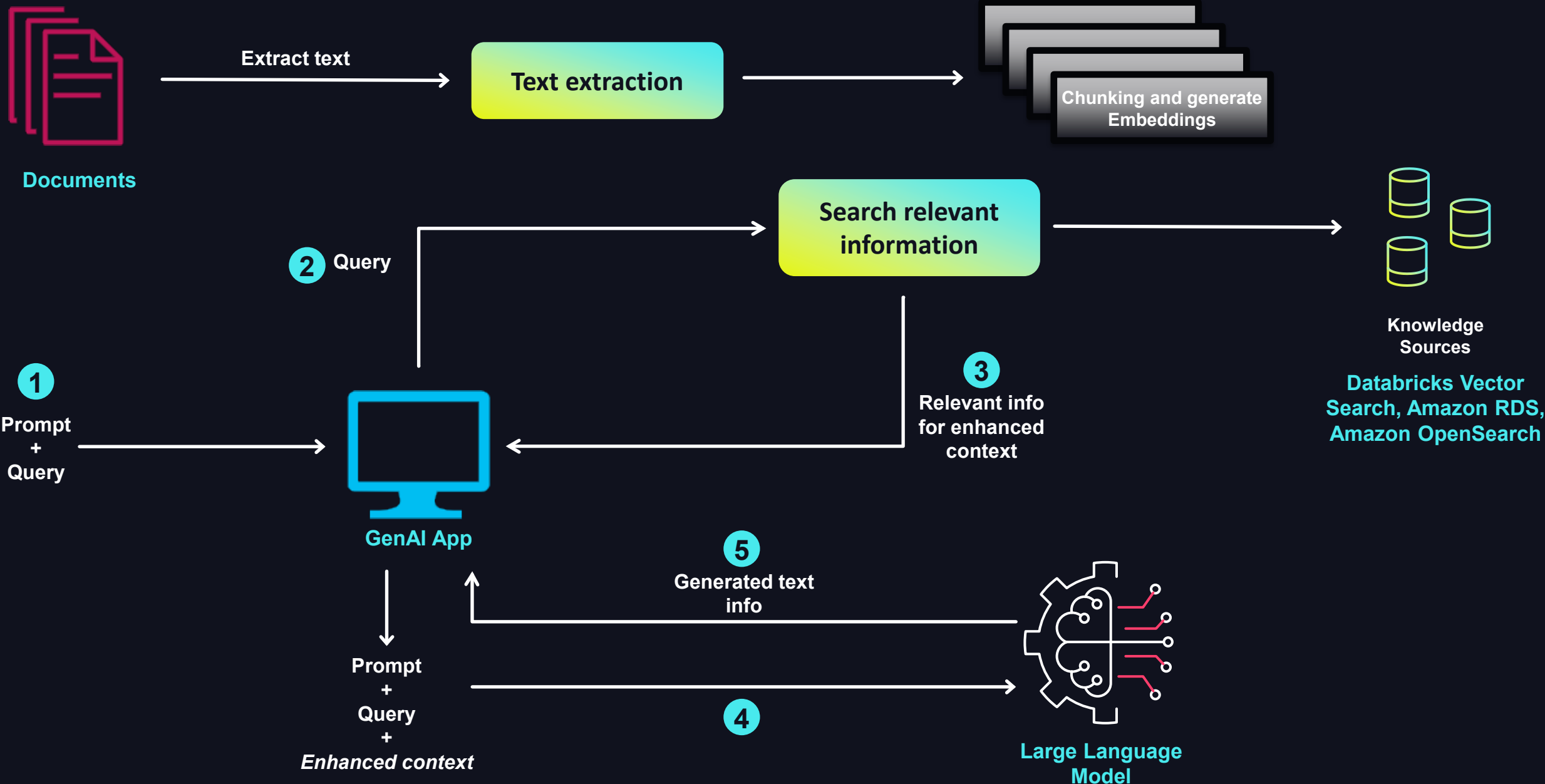
Question: What is AnyCompany's D2I ratio?
Text: Any company's 10K report
Answer: D2I ratio is 0.9%
Source: Page 22
===

Question: What are the company's financial risk factors?
Text: Any company's 10K report
Answer:

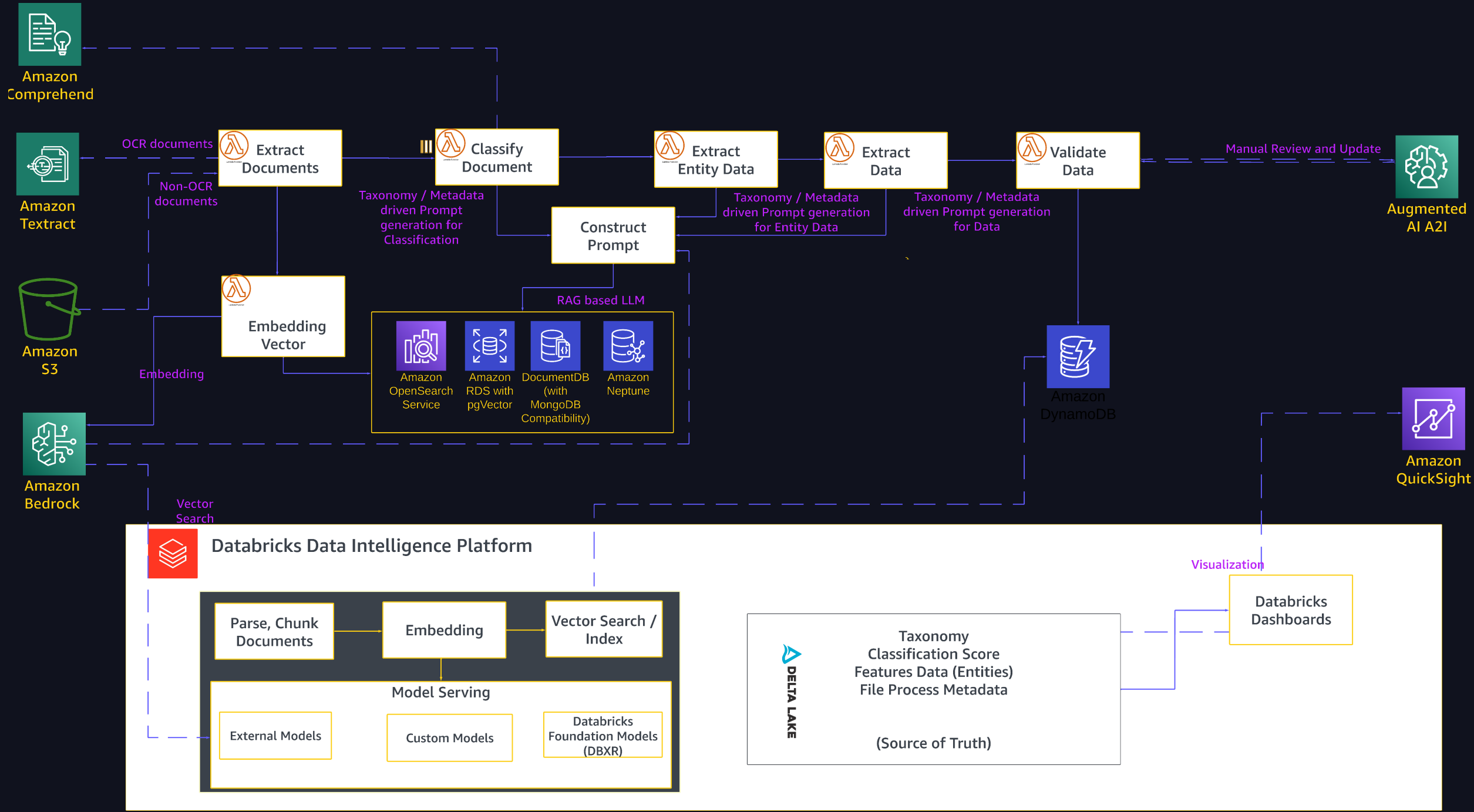


Generative AI - Using Retrieval Augmented Generation (RAG)

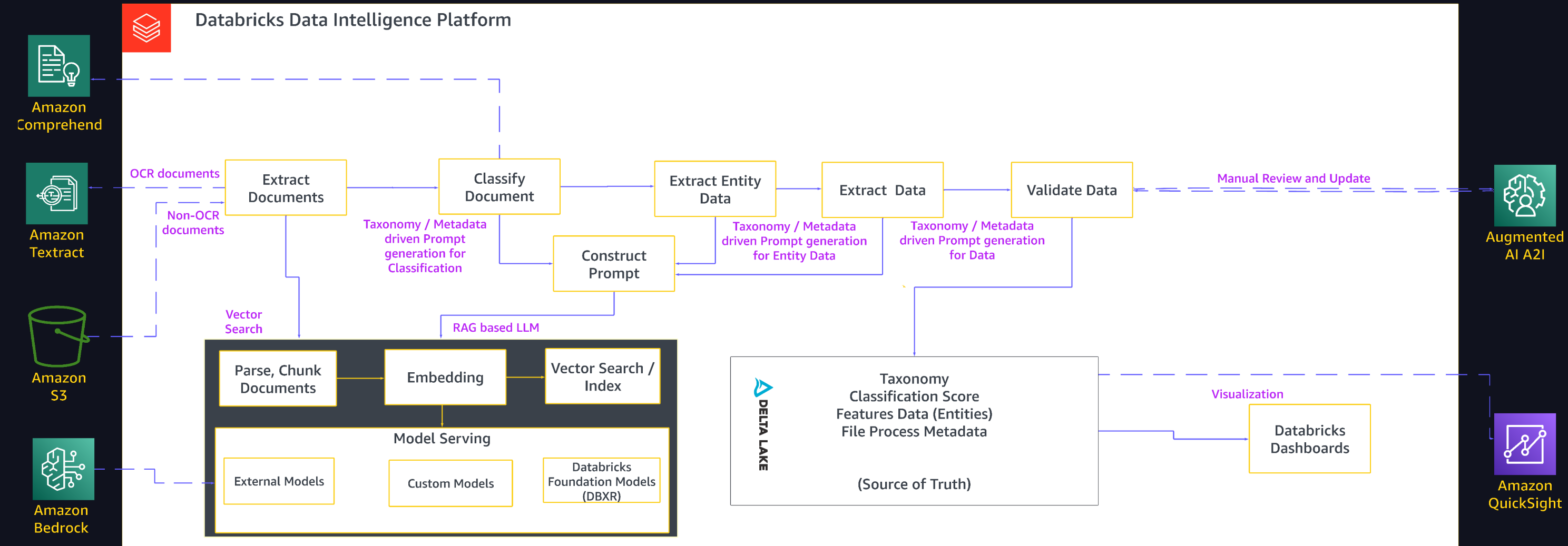
Augment prompts with relevant data in context



Reference Architecture – Using Databricks Approach (Batch, at-scale)



Reference Architecture – Using Databricks Approach (Batch, at-scale)



IDP Use-cases with Generative AI

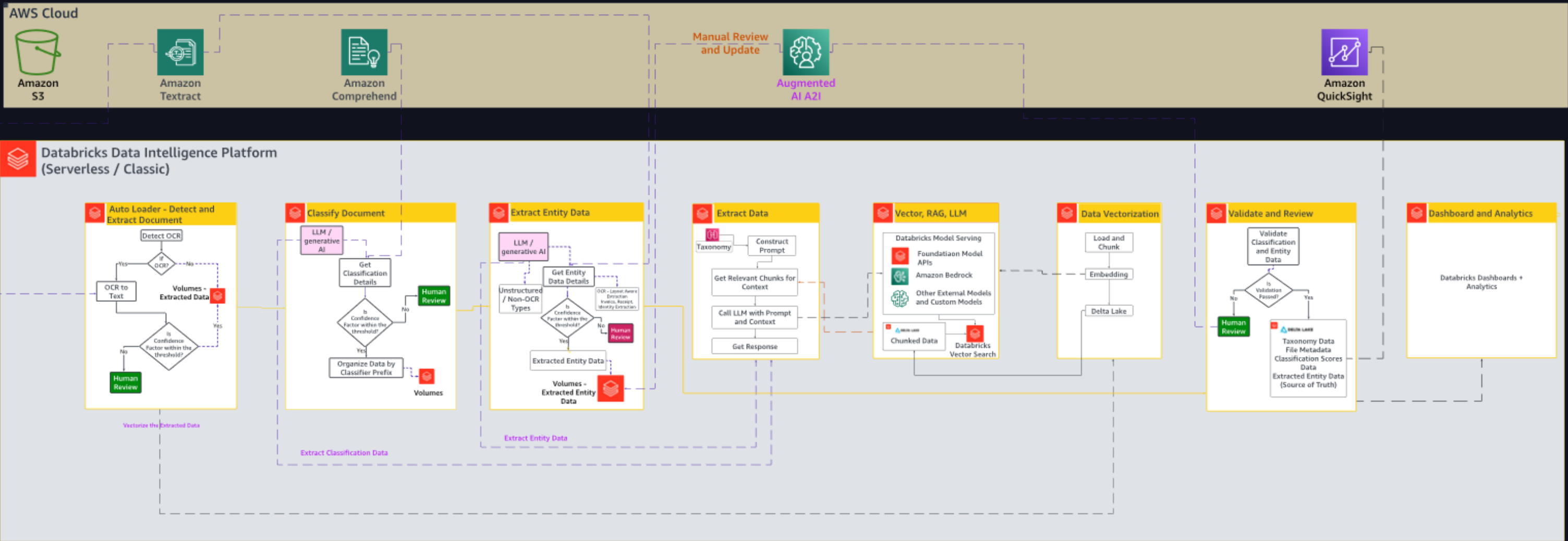
A few very common use cases

- Document Q&A with Chatbots
- Document summarization
- Enhanced data extraction
- Document classification
- Automated content creation
- Medical record analysis
- Translation and localization
- Learning and development

and more...



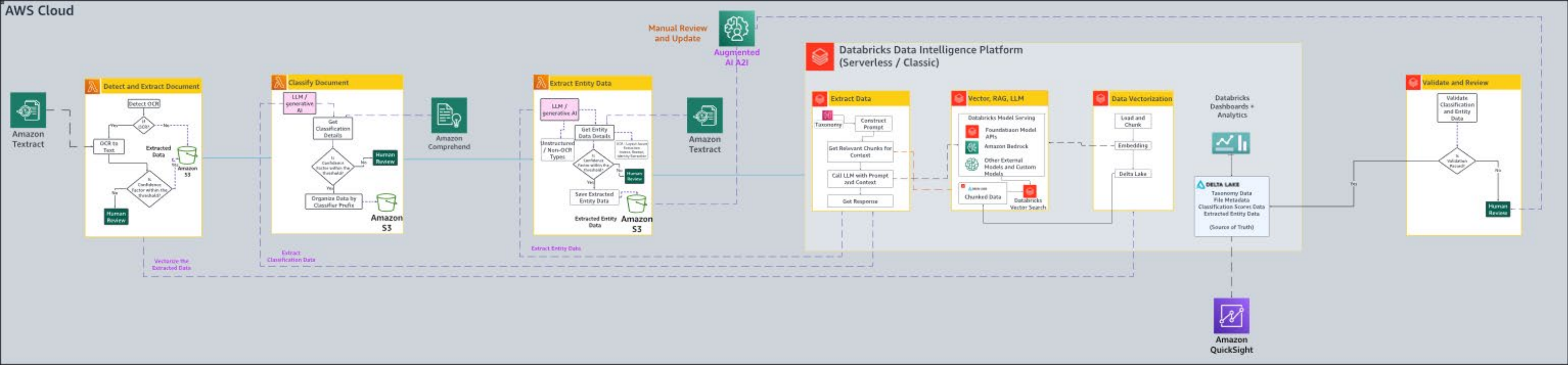
Reference Architecture – Using Databricks Approach (Batch, at-scale)



Delta Lake, Mosaic AI Vector Search, Deployment AI Gateway, DBRX and AWS AI Services and generative AI/LLM



Reference Architecture – Using AWS Lambda Approach (Real-time)



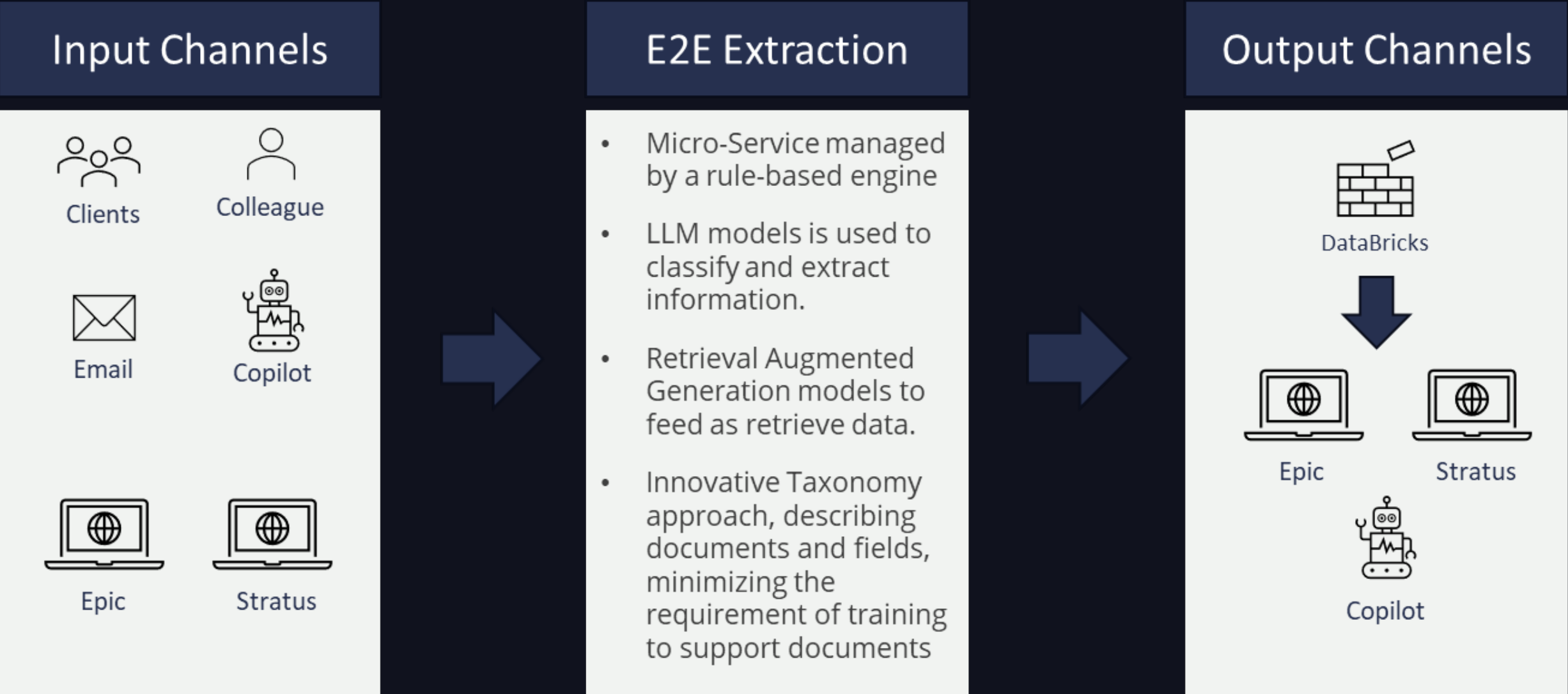
Amazon S3, AWS Lambda, SQS, SNS, Amazon Textract, Amazon Comprehend, Databricks' Delta Lake, Databricks Vector Store, Deployment AI Gateway, DBRX and 3P AI/LLM



Voice of Customer
Mayur Rajdev
Head of Architecture, AI & Automation
The Baldwin Group



End-to-End Enterprise Data Extraction



Taxonomy – generative AI augmented approach

Overview:

- The document Taxonomy is a representation of a document or group of documents that identifies and describes the document to be supported, it defines its characteristics, unique features to distinguish between other documents and entity fields to extract
- The taxonomy object is created using the Taxonomy Manager Ui, where the user can create and define the representation of a document
- This created JSON object will be used to create dynamic prompts using standard templates to feed the LLMs to classify and retrieve the relevant information from the document.

The screenshot shows a web interface for defining document taxonomy. It includes several sections:

- Form Fields:** Four dropdown menus for 'Division', 'Carrier', 'Document Type', and 'Policy Type', each set to 'Option 1'.
- Description:** A text area labeled 'Describe your document' containing the text: 'this is a loss run that contains XY this many pages, and these many elements'. Below it are 'Package' and 'Individual' radio buttons.
- Keywords:** A table titled 'Add unique elements that identify this document?' with columns 'Item', 'Description', and 'other'.

| Item | Description | other |
|-----------|--|-------------------------|
| Keyword 1 | you will find this in page 3 and will contain etc. | this is a table in text |
| Keyword 2 | you will find this in page 4 and will contain etc. | this is a table in text |
- Package documents:** A table with columns 'ID', 'Business', 'DocType', 'PolicyType', and 'Action'.

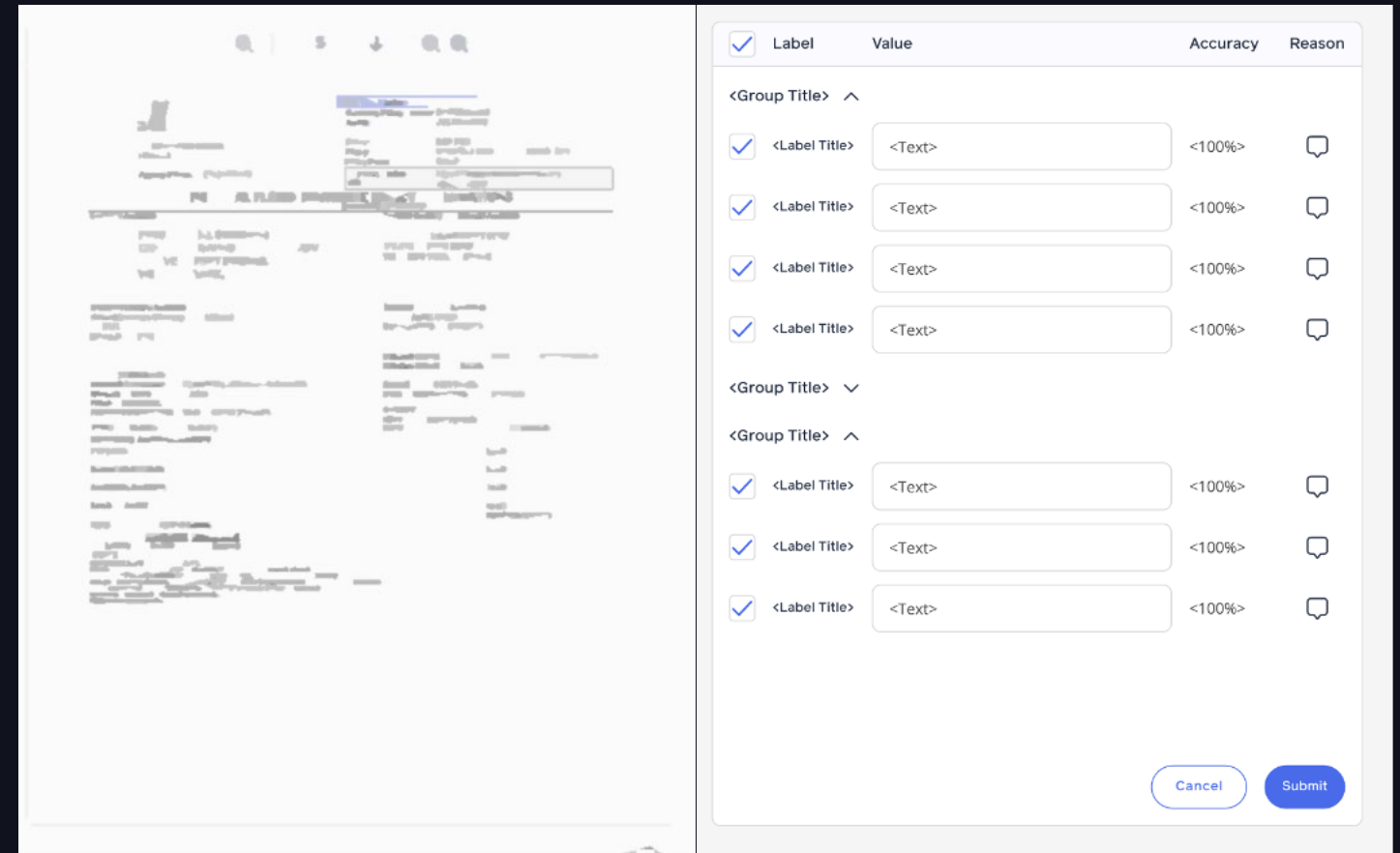
| ID | Business | DocType | PolicyType | Action |
|----|----------|---------|--------------|-----------------|
| 1 | CRMG | AUTO | Endorsement | {update delete} |
| 2 | CRMG | AUTO | Declarations | {update delete} |
- Extraction Fields:** A table with columns 'Item', 'Description', 'Alias', 'type', and 'Rule'.

| Item | Description | Alias | type | Rule |
|-----------------|---------------------------------------|----------|-----------|------------|
| Client | this is the client name found in blah | name etc | text | any rule |
| Expiration Date | this is the client name found in blah | name etc | date time | MM/DD/YYYY |

Classification and Extraction – generative AI augmented approach

Overview:

- The LLM will use the dynamic prompts generated by the taxonomy document, to classify the document into an individual file, or package as well as the type of document.
- Once classified, the system will use the defined entities from the taxonomy to individually search for each field in the document corpus and retrieve it in the proper format.
- Depending on the channel rules a Classification Station or Validation Station will be presented to a support team or to the business to confirm the accuracy of the classification or extraction and make the information available to the output channels and data lake.





Thank you!

Venkatavaradhan Viswanathan,
AWS, WW Databricks Architect
visvenky@amazon.com