

The background features a dark blue gradient with large, overlapping, semi-transparent shapes in shades of purple, pink, and orange, creating a modern, abstract design.

AWS re:Invent

NOV. 27 – DEC. 1, 2023 | LAS VEGAS, NV

AIM322-R

Intelligent document processing with generative AI for public sector

Sonali Sahu

Sr. Manager Specialist SA
AWS

Cameron Williams

Release Train Engineer
Booz | Allen | Hamilton



Agenda

Introduction

- Brief overview on intelligent document processing (IDP)

How to build large-scale document processing

- How Booz Allen built their document processing pipeline
- Large-scale document processing architecture

How generative AI can assist IDP

- Build document processing pipeline with generative AI

Wrap-up and next steps

Intelligent document processing (IDP)

STATE OF THE UNION



Accurate

Industry-leading accuracy among major on-premises and cloud vendors



Trustworthy

Elevate your security, compliance, and data privacy posture



Production ready

SLA-backed availability with high elasticity and scalability



Pay as you go

Pay for only what you use



Automate with confidence

[Paytm](#) extracts user data from documents with 97% accuracy



Improve business process efficiency

[Elevance Health](#) automated 90% of claims processing workflow



Lower costs

Customers achieved 73% ROI with AWS IDP ([Forrester Study](#))

Core IDP capabilities

IDP CAPABILITIES WITH AMAZON **TEXT**RACT



OCR



Layout



Forms



Entities



**Handwriting and
signatures**



Tables



**Queries/
custom
queries**



**Domain-specific
APIs (e.g., lending)**

Intelligent document processing demonstration



Need for an AI powered-IDP solution at VA

Use case

- 2022 PACT Act expanded benefits for over 5 million veterans
- The average benefits application packet contains 100+ documents averaging 13–15 pages each that can take ~30 minutes to review per document
- Expedite disability benefits determination by reducing the manual review of documents
- Extract raw text from the benefits application documents and enable a search engine that helps to quickly find and review the documents

Why Amazon IDP

- Available on GovCloud and is FedRAMP High
- Ability to scale as required to process 50–80 million pages per day
- High accuracy with both print and handwritten texts
- Reduce time spent reviewing data by highlighting results using Amazon Textract geometry data

Results/impact



- 400 million documents
- 10 billion+ pages processed
- Largest implementation of Amazon Textract in under one year



- Estimated 9x efficiency boost to document processing and manual review time



- Smart search engine to search and filter the content, minimizing number of documents to review

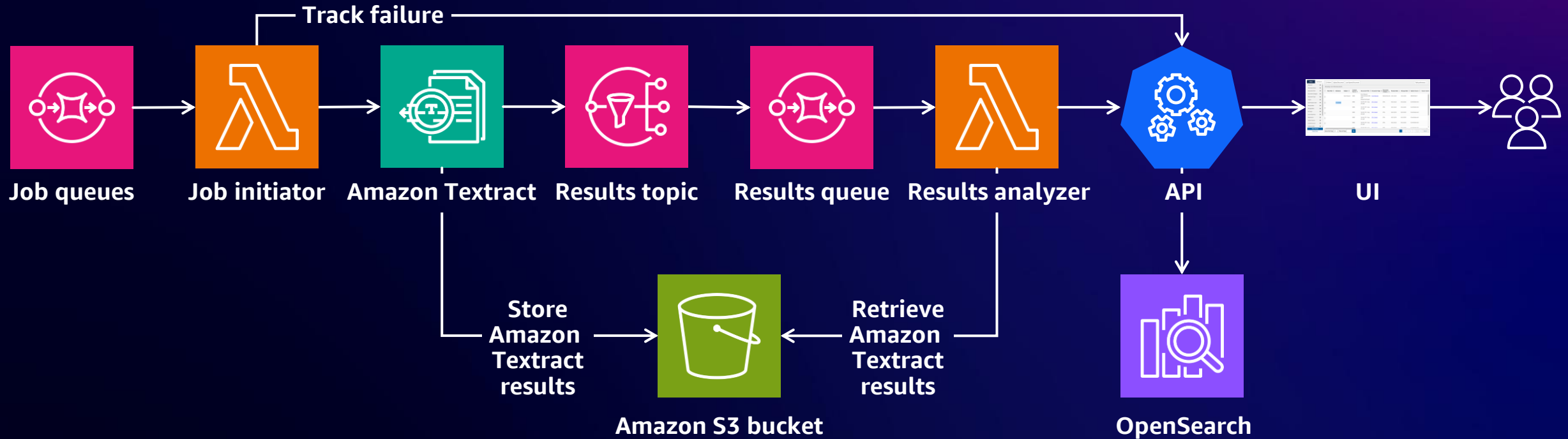


- Enabling future innovation through ML applications including generative AI by leveraging the extracted text

Questions

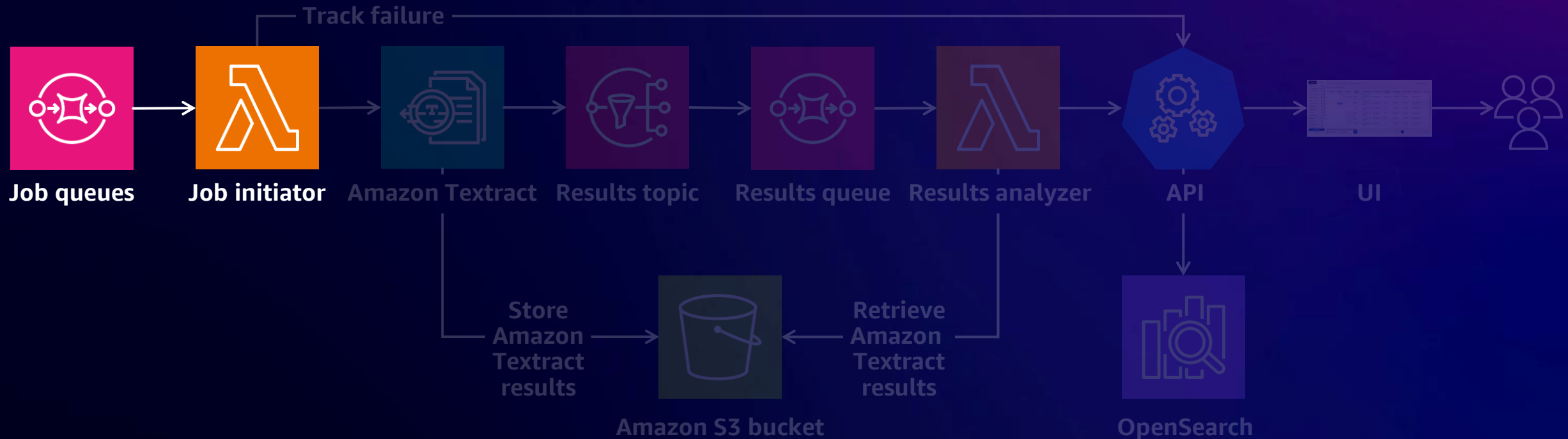


Our IDP workflow



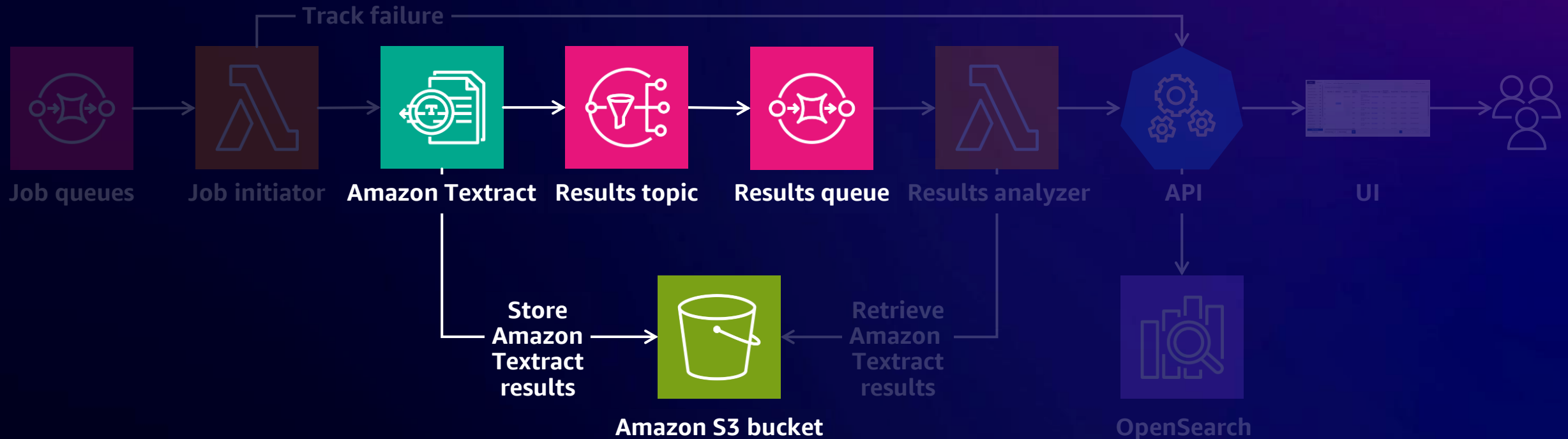
- System is FISMA High
- In GovCloud, all FedRAMP High
- Translation: services are secure and reliable

Our IDP workflow: Job queues and initiator



- Job queues invoke job initiator via event source mapping
- Amazon SQS and AWS Lambda are scalable, durable, and easy to integrate
- Consideration: balance message creation with rate limits

Our IDP workflow: Amazon Textract



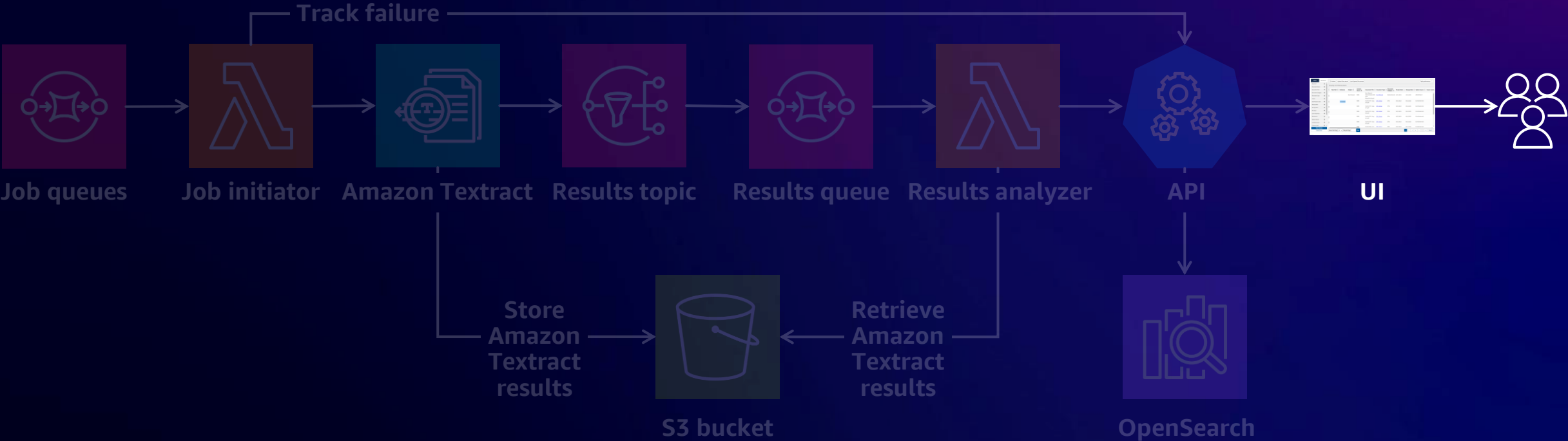
- Amazon Textract reads from and writes to S3 bucket using Lambda IAM role
- Synchronous (single-page) vs. asynchronous (multipage) operations
 - Synchronous operations return the result on the response, simplifying architecture
 - Asynchronous sends status to SNS (results topic)

Our IDP workflow: Results analyzer



- Results analyzer retrieves all pages of blocks from S3
- Maps results to custom internal schema and persists to RESTful API
- Store data in OpenSearch and S3

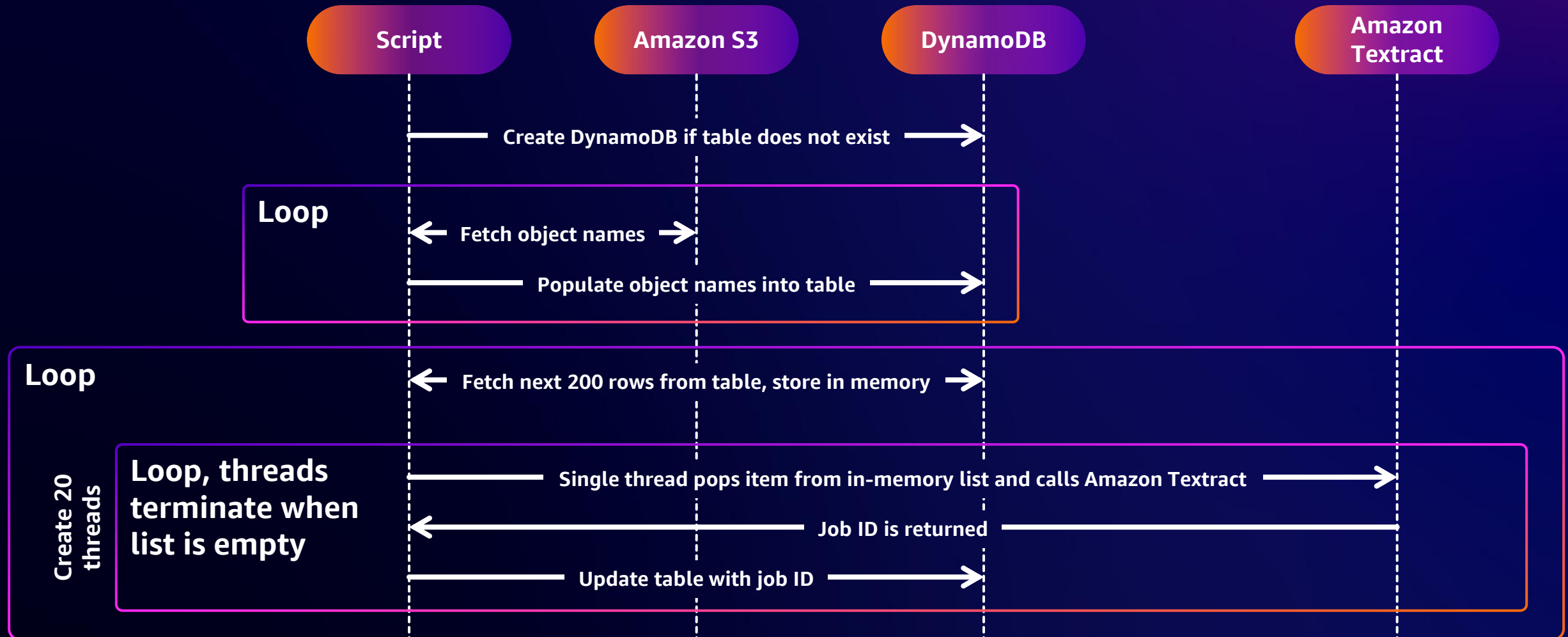
Our IDP workflow: User interaction



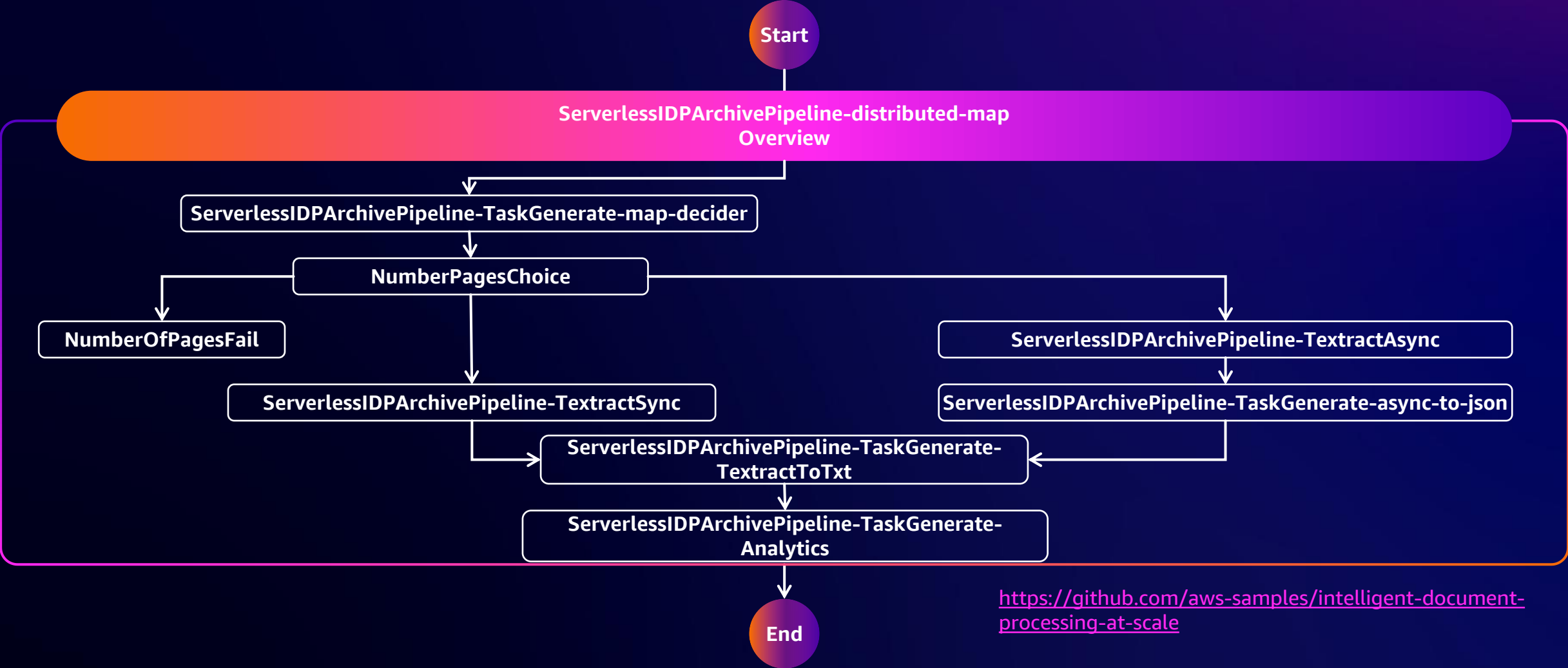
Demonstration



Large-scale design flow



Large-scale serverless deployment – IDP



<https://github.com/aws-samples/intelligent-document-processing-at-scale>



Questions

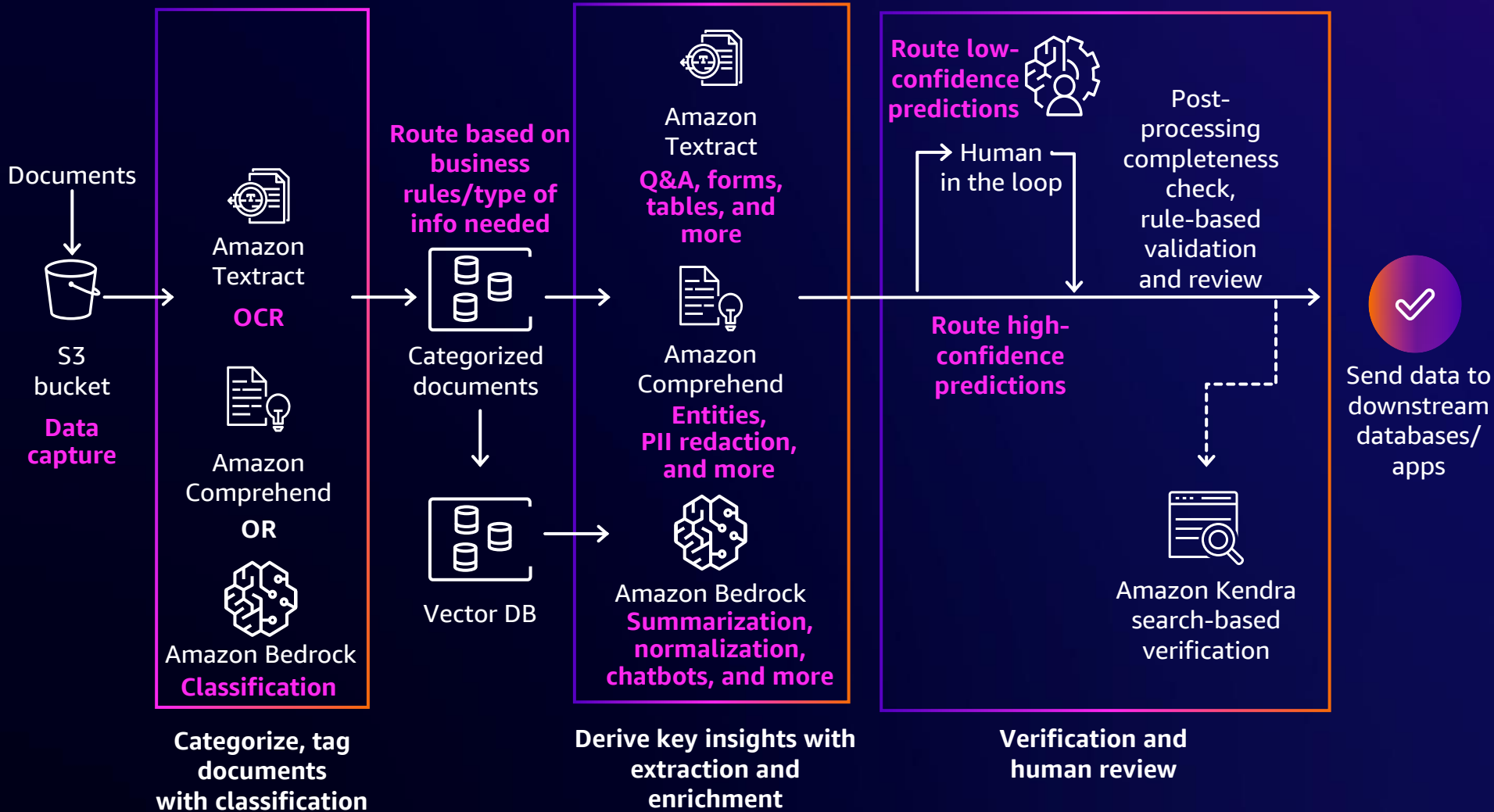


Moving forward with generative AI

- Built prototype to answer predefined questions
- Collaborating with AWS on personally identifiable information recognition models as a step toward keeping generative AI solutions compliant
- Future possibilities
 - Prepopulate long, complex forms based on model findings
 - Generate recommended adjudication of claims
 - Racial and gender bias reports

AI is here; let's generate a better world!

Document pipeline with AWS IDP and generative AI



How to get started

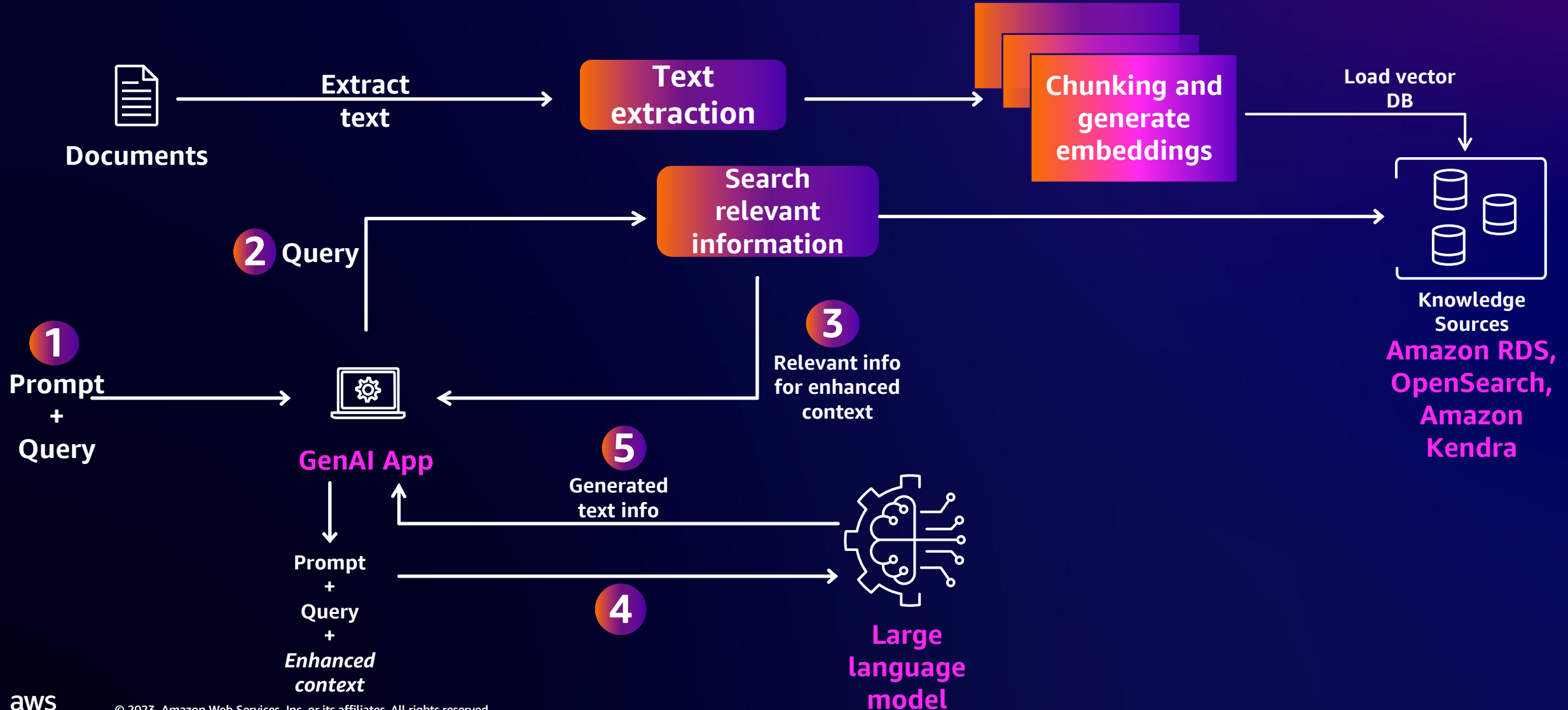
Developers can easily embed AI-powered functionality from Textract and Comprehend into your business workflows and apps

Engage your data science team for FM selection, evaluation and tuning based on your generative AI use case

Link your gen AI/FM modules with AWS IDP (e.g., through chaining) to create an end-to-end document processing pipeline

Using retrieval-augmented generation (RAG)

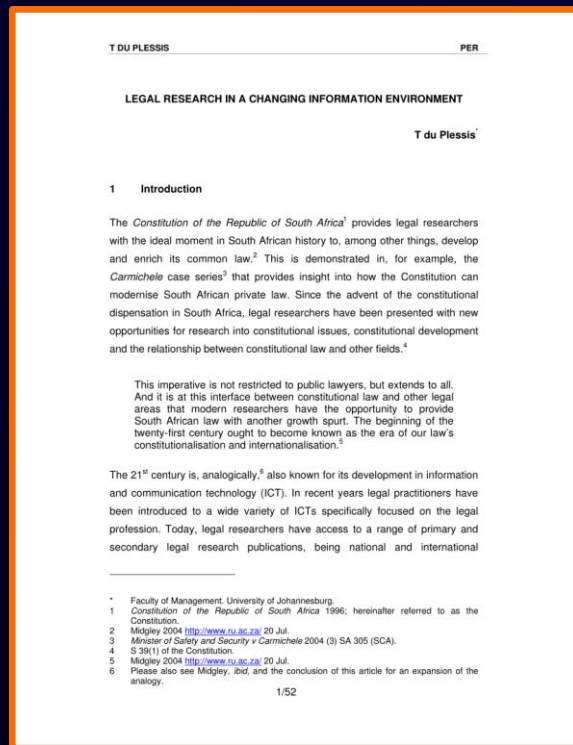
AUGMENT PROMPTS WITH RELEVANT DATA IN CONTEXT



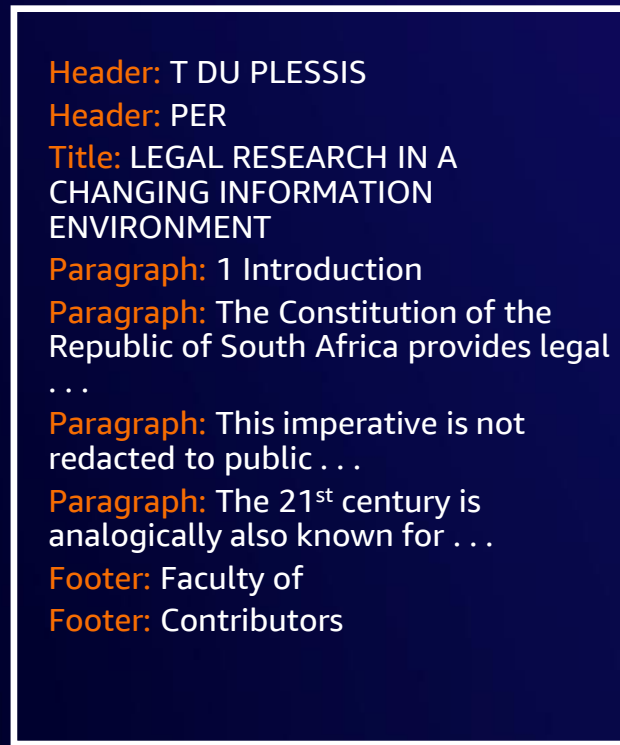
Layout detection for improved LLM accuracy

LAUNCHED!

DOCUMENT



OUTPUT



Extract titles, paragraphs, footer, headers, etc.



Orders output in a way humans would read (reading order)



Reduced post-processing when extracting information from dense documents

Example

LAYOUT_HEADER

RecSys '23, September 18-22, 2023, Singapore, Singapore

Elad Haramaty, Zohar Karnin, Arnon Lazerson, Liane Lewin-Eytan, and Yoelle Maarek

LAYOUT_TITLE

Extended Conversion: Capturing Successful Interactions in Voice Shopping

5 ECVR VS CVR

LAYOUT_SECTION_HEADER

After setting the parameters of the ECVR metric, we provide a deeper comparison between ECVR and the standard immediate conversion metric (CVR). We demonstrate ECVR superiority in terms of sensitivity and long-term effect. In addition, we show that a ranker optimized for ECVR outperforms a ranker optimized for CVR.

Sensitivity. We repeat the experiment pertaining to sensitivity that was described in Section 4.2, this time measuring ECVR and CVR induced by applying random and relevance-based ranking to voice shopping traffic corresponding to early phases of the shopping journey. We expect a sensitive metric to give a higher score to the relevance-based ranker as it provides better user experience. Indeed, as both metrics are sensitive, they give a higher score to the relevance-based ranker. The score difference between the relevance ranker and the random ranker is $0.51 \pm 0.15\%$ for ECVR and $0.09 \pm 0.01\%$ for CVR with a CI of 95%. As evident by the large and statistically significant difference, the ECVR metric is more sensitive to changes in user experience.

LAYOUT_SECTION_TEXT

6 DEFINING EXTENDED CONVERSION AS A PARAMETERIZED METRIC

We consider a hierarchy of five natural product similarity levels differing in their specificity, from the most specific similarity level, where a product is only similar to itself, to the most general similarity level, which includes all products:

- Product: a trivial similarity level in which a product is only similar to itself.
- Substitutions: products are considered similar if they can replace one another. In other words, the customers are mostly indifferent between substituting products.
- Type: products are considered similar if they belong to the same product-type.
- Department: products are considered similar if they belong to the same department, which is a natural way to partition the universe of products, just like aisle descriptions in a physical store.
- All: a trivial similarity level, in which a product is similar to all other products.

LAYOUT_LIST

RecSys '23, September 18-22, 2023, Singapore, Singapore

Elad Haramaty, Zohar Karnin, Arnon Lazerson, Liane Lewin-Eytan, and Yoelle Maarek

Extended Conversion: Capturing Successful Interactions in Voice Shopping

5 ECVR VS CVR

After setting the parameters of the ECVR metric, we provide a deeper comparison between ECVR and the standard immediate conversion metric (CVR). We demonstrate ECVR superiority in terms of sensitivity and long-term effect. In addition, we show that a ranker optimized for ECVR outperforms a ranker optimized for CVR.

Sensitivity. We repeat the experiment pertaining to sensitivity that was described in Section 4.2, this time measuring ECVR and CVR induced by applying random and relevance-based ranking to voice shopping traffic corresponding to early phases of the shopping journey. We expect a sensitive metric to give a higher score to the relevance-based ranker as it provides better user experience. Indeed, as both metrics are sensitive, they give a higher score to the relevance-based ranker. The score difference between the relevance ranker and the random ranker is $0.51 \pm 0.15\%$ for ECVR and $0.09 \pm 0.01\%$ for CVR with a CI of 95%. As evident by the large and statistically significant difference, the ECVR metric is more sensitive to changes in user experience.

6 DEFINING EXTENDED CONVERSION AS A PARAMETERIZED METRIC

We consider a hierarchy of five natural product similarity levels differing in their specificity, from the most specific similarity level, where a product is only similar to itself, to the most general similarity level, which includes all products:

.....

***Extracted text trimmed for brevity*

<https://aws.amazon.com/blogs/machine-learning/amazon-textracts-new-layout-feature-introduces-efficiencies-in-general-purpose-and-generative-ai-document-processing-tasks/>



Questions



Getting started



**Get started with
Amazon Bedrock**



**IDP hands-on
Workshop
with Layout**



**Get started with
hands-on
Amazon Bedrock
workshop**



**Get started with
hands-on
IDP OpenSearch
workshop**

Thank you!

Sonali Sahu

sonalsah@amazon.com

Cameron Williams



Please complete the session survey in the mobile app