# Amazon Textract

## Specify and extract information using the Queries feature in Amazon Textract

Navneeth Nair

Senior Product Manager, Amazon Textract
AWS

Martin Schade

Senior ML Product SA, Amazon Textract
AWS

aws

# Agenda

Overview of AI/ML at AWS

Issues with legacy solutions for document processing
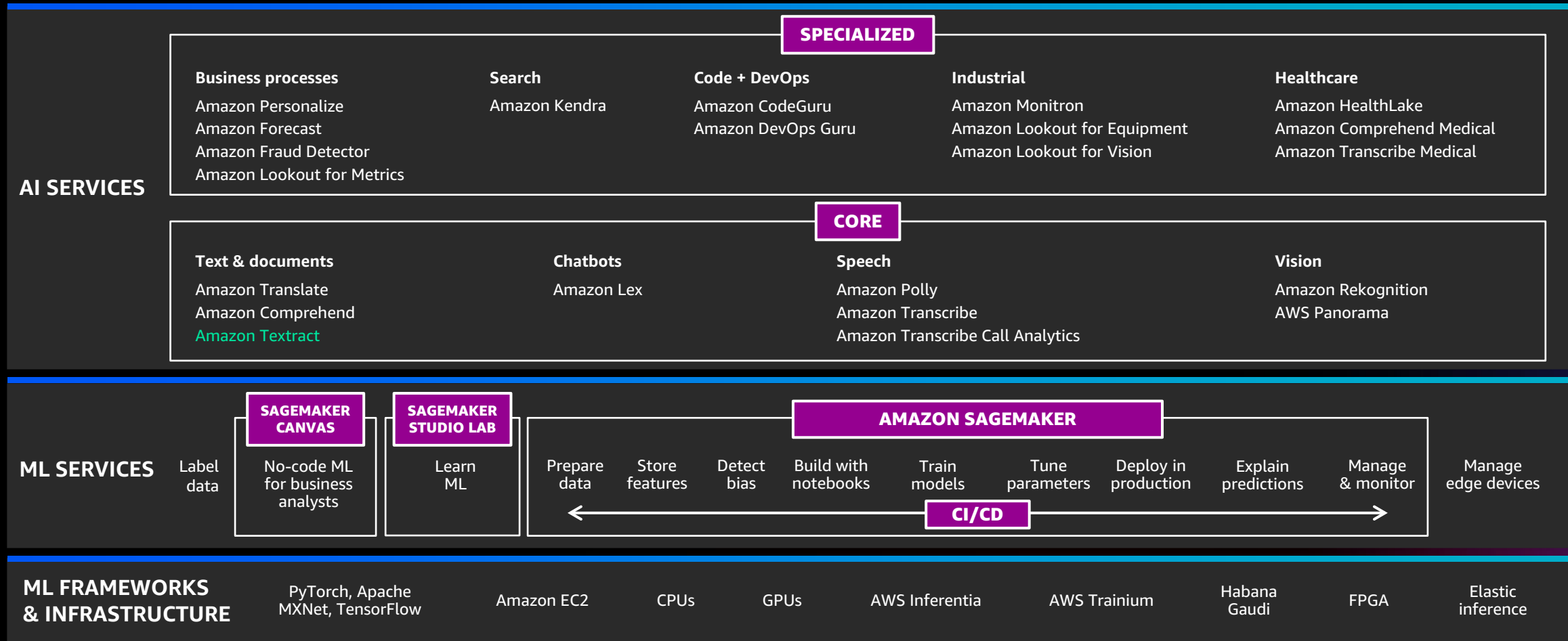
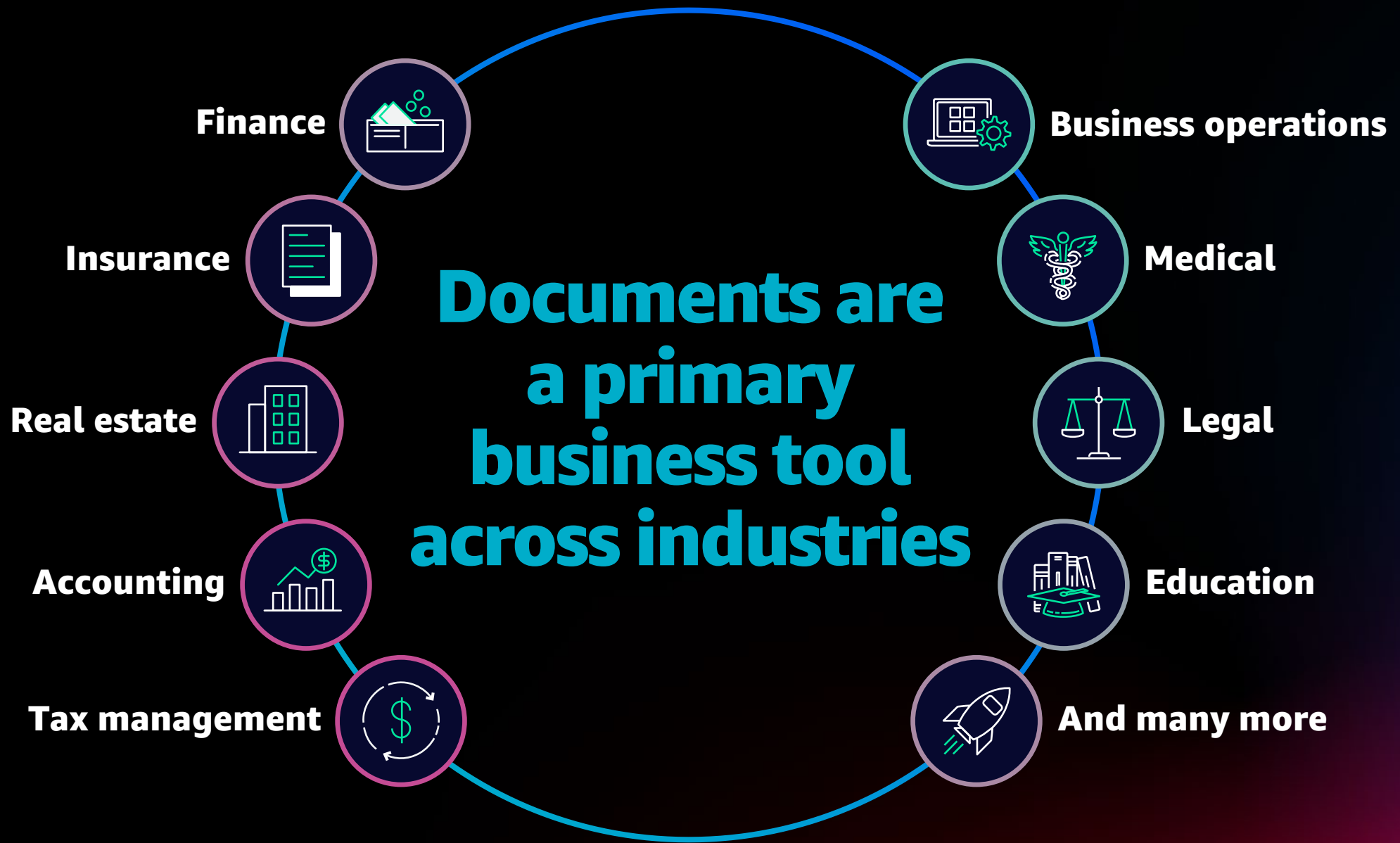Amazon Textract's current capabilities

Textract Analyze Document – "Queries"

Demo: Examples of automated data extraction from mortgage and insurance industry documents
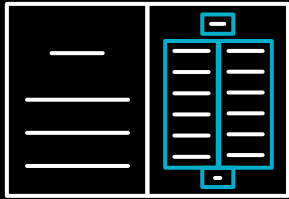
# The AWS AI/ML stack

## BROADEST AND MOST COMPLETE SET OF MACHINE LEARNING CAPABILITIES

**AI SERVICES**

**SPECIALIZED**

| **Business processes** | **Search** | **Code + DevOps** | **Industrial** | **Healthcare** |
|---|---|---|---|---|
| Amazon Personalize | Amazon Kendra | Amazon CodeGuru | Amazon Monitron | Amazon HealthLake |
| Amazon Forecast | | Amazon DevOps Guru | Amazon Lookout for Equipment | Amazon Comprehend Medical |
| Amazon Fraud Detector | | | Amazon Lookout for Vision | Amazon Transcribe Medical |
| Amazon Lookout for Metrics | | | | |

**CORE**

| **Text & documents** | **Chatbots** | **Speech** | **Vision** |
|---|---|---|---|
| Amazon Translate | Amazon Lex | Amazon Polly | Amazon Rekognition |
| Amazon Comprehend | | Amazon Transcribe | AWS Panorama |
| Amazon Textract | | Amazon Transcribe Call Analytics | |

**ML SERVICES**

**SAGEMAKER CANVAS**

**SAGEMAKER STUDIO LAB**

**AMAZON SAGEMAKER**

| Label data | No-code ML for business analysts | Learn ML | Prepare data | Store features | Detect bias | Build with notebooks | Train models | Tune parameters | Deploy in production | Explain predictions | Manage & monitor | Manage edge devices |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

**CI/CD**

**ML FRAMEWORKS & INFRASTRUCTURE**

| PyTorch, Apache MXNet, TensorFlow | Amazon EC2 | CPUs | GPUs | AWS Inferentia | AWS Trainium | Habana Gaudi | FPGA | Elastic inference |
|---|---|---|---|---|---|---|---|---|

aws

Documents are
a primary
business tool
across industries

Finance

Insurance

Real estate

Accounting

Tax management

Business operations

Medical

Legal

Education

And many more

# Legacy document processes do not meet today's needs

Legacy OCR and manual processes are time-consuming, error-prone, and expensive

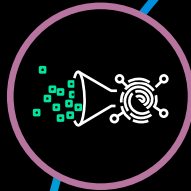Manual processes do not scale easily with document volume

Difficult to find useful information needed for business decisions

**Faster document processing shortens decision cycles so you can serve more customers and have people do higher-value tasks**

# Amazon Textract

**Go beyond OCR with accurate, versatile information extraction**

**No ML experience required**

**Reduce manual effort**

**Lower costs**

## Amazon Textract

Automatically extract printed text, handwriting, and data from any document

# Key verticals and business outcomes

### Financial services

Accurately extract data and insights from pay stubs, loan applications, tax forms, driver's licenses, invoices, titles, etc.

### Manufacturing and retail

Analyze invoices, bills of materials, contracts, licenses, warranties, and other agreements

### Healthcare and life sciences

Process insurance claims and quickly get to clinical insights faster from healthcare and patient intake forms

### Public sector

Improve speed and quality of services by gathering data from healthcare records, government forms, public records, taxes, securities filings, and unstructured documents

### Productivity

Improve accuracy and process velocity and reduce costs

### End-user experience

Automate user data capture and provide a better end-user experience through convenience, speed of service, and improved availability

### Decision-making

Improve decision-making through richer inferences to scale business processes more effectively

# Improving loan underwriting process productivity

### Problem
Incomplete loan packages, missing data, and the discovery of additional documentation requirements during the underwriting process often creates more work and slows the process down, resulting in fewer loans reviewed by underwriters and added lender costs

### Solution
Black Knight developed Underwriter Assist (UA), an AI-powered mortgage solution that leverages Amazon Textract and other machine learning technologies; UA performs document recognition, data extraction, and analysis of income, assets, and property appraisal

### Impact
Using advanced machine learning services like Amazon Textract, Black Knight has increased lender productivity, driven operational efficiency, and reduced costs, resulting in higher customer satisfaction

**BLACK KNIGHT**

# Scaling claims processing through automation

### Problem
Health insurance companies spend millions of dollars to extract sensitive information from claims forms and accompanying attachments to perform their business operations; to reduce manual labor, the company wanted to automate the process

### Solution
Anthem chose Amazon Textract to digitize and automate its claims process, for its image-processing capability, ability to detect tables and forms, and adherence to security and compliance standards
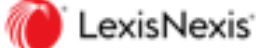
### Impact
Using Amazon Textract, Anthem can extract and digitize data to quickly process thousands of claims each day; having already automated 80% of their claim processing workflow, Anthem expects to reach 90% automation levels or higher on AWS

# Amazon Textract customers

| | | | | | |
|---|---|---|---|---|---|
| 3M | Anthem | ARQ | Assent | AXRAIL | bakertilly |
| BLACK KNIGHT | BDO | beyond lucid technologies | BlueVine | Broadridge | Canara HSBC OBC Life Insurance |
| CHANGE HEALTHCARE | CHISEL | Cox AUTOMOTIVE | THE GLOBE AND MAIL* | healthfirst | HELLOSIGN |
| HIVEO | IHS Markit | intuit | LUMIQ | Met Office | MOODY'S ANALYTICS |
| NHS | PennyMac | pfs TECH | QL Resources Berhad | Roche | rekeep |
| salesforce | TC Energy | Wrapped Insurance | FRED HUTCH | LexisNexis | FINRA |

# Amazon Textract:
## Generally available features overview

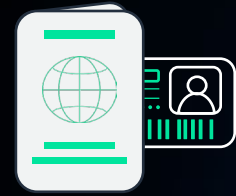# Amazon Textract capabilities



**Text**



**Forms**

Invoices and receipts

Identity documents

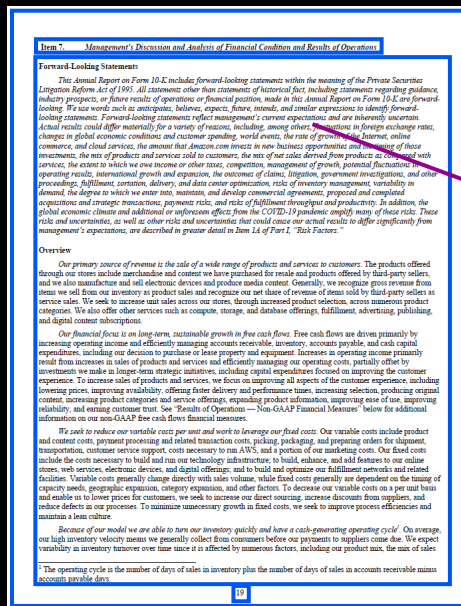

**Specialized documents**



**Handwriting**



**Tables**

Earnings

Paystubs



**Queries**

# Text extraction

## OPTIMAL FOR DENSE TEXT EXTRACTION WITH INDUSTRY-LEADING OCR ACCURACY

**Document**

**Output**

Blocks

Page, line, word

is washed by waves, and cooled

Word    Line 1

Outputs detected text in 3 hierarchy blocks: page, line, and word

Bounding box for each line and word provide visual cues for post-processing

Included confidence scores enable informed decision-making for your workflows
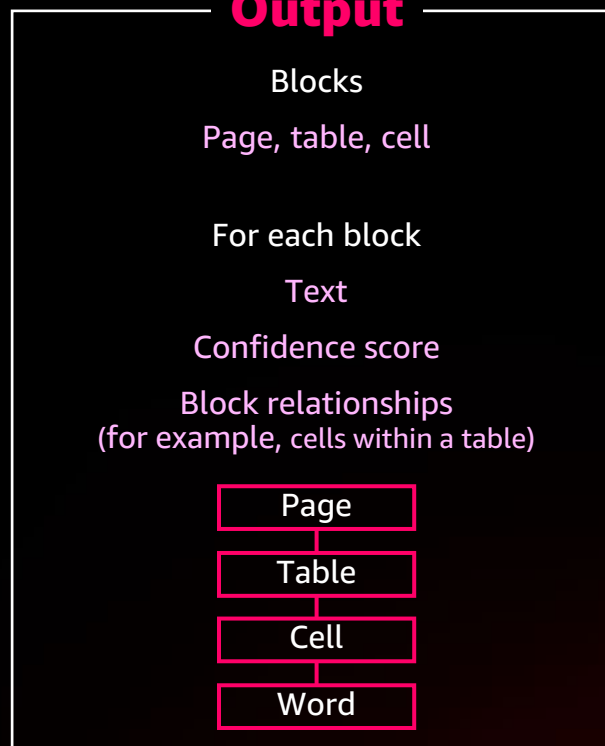
# Table extraction

EXTRACT TABLES FROM DOCUMENTS WHILE PRESERVING DATA STRUCTURE AND RELATIONSHIPS

## Document

| Previous employment history | | | | |
|---|---|---|---|---|
| Start date | End date | Employer name | Position held | Reason for leaving |
| 1/15/2009 | 6/30/2013 | Any company | Head Baker | Family relocated |
| 8/15/2013 | Present | Example corp. | Baker | N/A, current employer |

## Output

Blocks

Page, table, cell

For each block

Text

Confidence score

Block relationships
(for example, cells within a table)

Page

Table

Cell

Word

✓ Outputs recognized tables with relationships data intact

✓ Intelligently groups cells within tables and words within each cell

✓ Output also includes confidence scores, geometry info, and row/column indexes

# Form extraction

## Document

Full Name

| John | X | Doe |
|---|---|---|
| First | Middle | Last |

Date of Birth

| 01 | 01 | 1971 |
|---|---|---|
| MM | DD | YYYY |

Gender
Male ●
Female ○

## Output

Blocks

Page, key-value set

Example output

First: John

Middle: X

Last: Doe

MM: 01
DD: 01
YYYY: 1971

Male: True
Female: False

Outputs form field name (key) and field value name (value) with relationship data intact

Captures logical groupings, relationships, and glyphs

Output also includes confidence scores and geometry info

# Invoices and receipts

SPECIALIZED SUPPORT TO PROCESS INVOICES AND RECEIPTS AT SCALE

## Document



→

## Output

Summary fields
Vendor Name: WHOLE FOODS MARKET
Subtotal (SUBTOTAL): $3.50
Net Sales (OTHER): $3.50
Tax/Fee (TAX): $0.39
Sold Items (OTHER): 1
Paid (OTHER):
Debit (OTHER): $3.89
Tax/Fee Total (TAX): $0.39
Total (TOTAL): $3.89

Line items
ITEM: Pizza Slice
PRICE: $3.50

✓ Outputs headline amounts, line item details, and inferred fields (like Vendor Name)

✓ Supports any style of invoice or receipt

✓ No templates or configuration required

# Identity documents

SPECIALIZED SUPPORT FOR IDENTITY DOCUMENTS

## Document



## Output

First Name: JORGE
Last Name: SOUZA
Middle Name:
Address Line1: 100 MAIN STREET
Address Line 2:
City: ANYTOWN
State: MA
Document Number: 820BAC729CBAC
Expiration Date: 01/20/2020
Date of Birth: 03/18/1978
ID Type: Driver License
Date of Issue: 03/18/1978
Issued By: MASSACHUSETTS
Class: D
Restrictions: NONE
Endorsements: NONE

✓ Over 95% accuracy for US driver licenses and passports

✓ No templates or configuration required

✓ Outputs normalized field names and supports implied elements

# Textract Queries

# Data extraction challenges

## Data Variations

**(SSN vs. Social Security Number vs. Tax ID)**

Post processing

Expensive

Time consuming

## Data Structure

**(Table, Form, Implied Fields)**

Post processing

Custom models

Manual

## Nested Data

**(page sections, duplicate fields)**

Development and management overhead

Templates are brittle

# Textract queries

EASILY SPECIFY AND EXTRACT VALUABLE PIECES OF INFORMATION FROM DOCUMENTS

## Document

## Questions

Sample questions
1. What is the patient's first name?
2. When was the first dose administered?
3. What are YTD earnings in this pay stub?

and more

## Output

1. What is the patient's first name? [Patient first name]
2. When was the first dose administered? [Date]
3. What are YTD earnings in this pay stub? [YTD earnings]

…

*No custom model training or synonym lists required*

Combines visual, spatial, and textual cues to provide higher accuracy

Simple Q&A response reduces post-processing

Extract only the specific fields that matter

Easy integration with existing Amazon Textract Analyze Document API

# Example: Amazon Textract Queries on pay stubs

**Pay stubs**



**Example – tabular data extraction**
- What is the gross pay this period? $452.43
- What is the gross pay YTD? $23,526.80
- What is the Medicare tax YTD? $341.12

**Example – implied fields**
- What is the company name? ANY COMPANY CORP.
- What is the employee name? JOHN STILES

# Example: Textract Queries on Fannie Mae Form 1003



Uniform Residential Loan Application (Form 1003)

**For example**
- What is the borrower's date of birth? **01/01/1900**
- What is the co-borrower's date of birth? **02/02/1902**

# Example: Textract Queries on health insurance cards

**Health insurance card**



**For example**
- What is the member name? **Jacob Michael**
- What is the plan type? **AnyPlan X-EPO**
- What is the insurance provider? **AnyInsuranc Co.**

# Demo

# Resources

## Amazon Textract service

- [https://aws.amazon.com/textract/](https://aws.amazon.com/textract/)

## Blogs

- Automatically extract text and structured data from documents with Amazon Textract – [https://go.aws/3mmDXUI](https://go.aws/3mmDXUI)
- Specify and extract information from documents using the new Queries feature in Amazon Textract – [https://go.aws/3xpBShn](https://go.aws/3xpBShn)
- Announcing support for extracting data from identity documents using Amazon Textract – [https://go.aws/3Q1eWvX](https://go.aws/3Q1eWvX)
- Announcing specialized support for extracting data from invoices and receipts using Amazon Textract – [https://go.aws/3miPzln](https://go.aws/3miPzln)

## GitHub repository

- Mortgage – [https://bit.ly/3thI7Bs](https://bit.ly/3thI7Bs)
- Insurance – [https://bit.ly/3mqEV2d](https://bit.ly/3mqEV2d)
- Vaccination cards – [https://bit.ly/3aJJlPo](https://bit.ly/3aJJlPo)
- Constructs Hub: [https://constructs.dev/packages/amazon-textract-idp-cdk-constructs/](https://constructs.dev/packages/amazon-textract-idp-cdk-constructs/)
- GitHub Constructs: [https://github.com/aws-samples/amazon-textract-idp-cdk-constructs/](https://github.com/aws-samples/amazon-textract-idp-cdk-constructs/)
- GitHub Samples: [https://github.com/aws-samples/amazon-textract-idp-cdk-stack-samples](https://github.com/aws-samples/amazon-textract-idp-cdk-stack-samples)

# Thank you!