



# AI Driven Data Enrichment with Azure Cognitive Search

---

Matt How | Adatis



# Matt How

Head of Data Science

- 8 Years working in Data & AI on the Microsoft Data Platform
- Broad experience working with Azure Cognitive Search and the Applied AI Components
- Published author on the topic of the Modern Data Warehouse in Azure



# Sponsors – Thank you



Part of Telefónica Tech



**ADVANCING  
ANALYTICS**



**PRIMUS CONNECT**

*We couldn't do it without you.*



# Agenda

---

- **Cognitive Search Building Blocks**
- **Custom Skills & AI Enrichment**
  - **Cognitive Skills**
  - **Custom Skills**
  - **Azure OpenAI & Document Intelligence Integration**
  - **Vector Embeddings**
- **Enhancing Cognitive Search Usage**

# Cognitive Search Building Blocks

# What is Cognitive Search?



## Enterprise Search

- Instant document retrieval
- Intuitive search features
- Accurate search results
- Optimised search algorithms



## Knowledge Mining

- Surface essential data points
- Semantic search capability
- Intelligent document ranking
- Optical Character Recognition

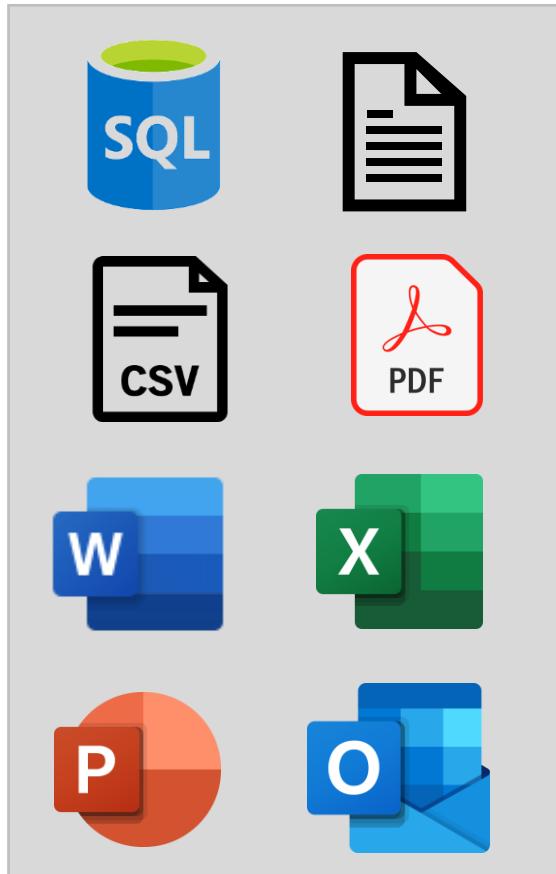


## Document Intelligence

- Rich document metadata
- AI driven document enrichment
- Document classification
- Digitising assets for other uses

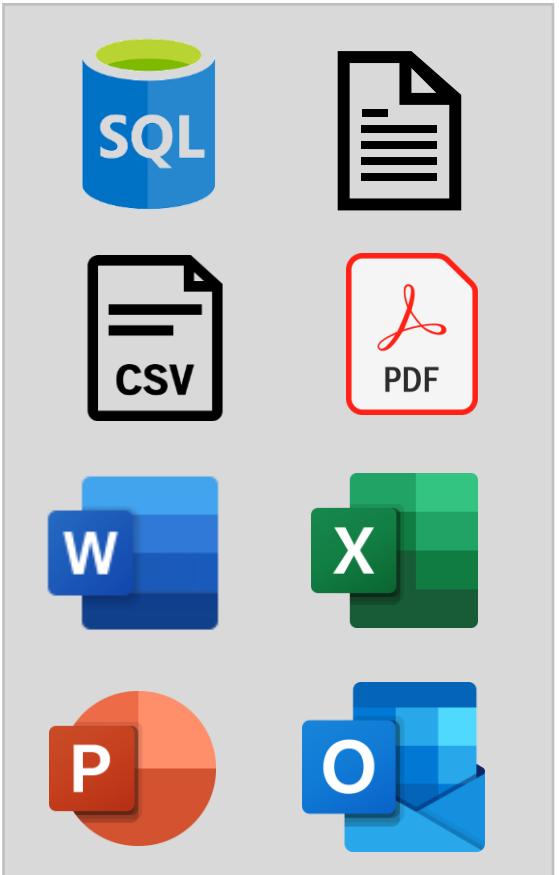
# Cognitive Search Building Blocks – Data Source

## Data Source



# Cognitive Search Building Blocks – Index

## Data Source

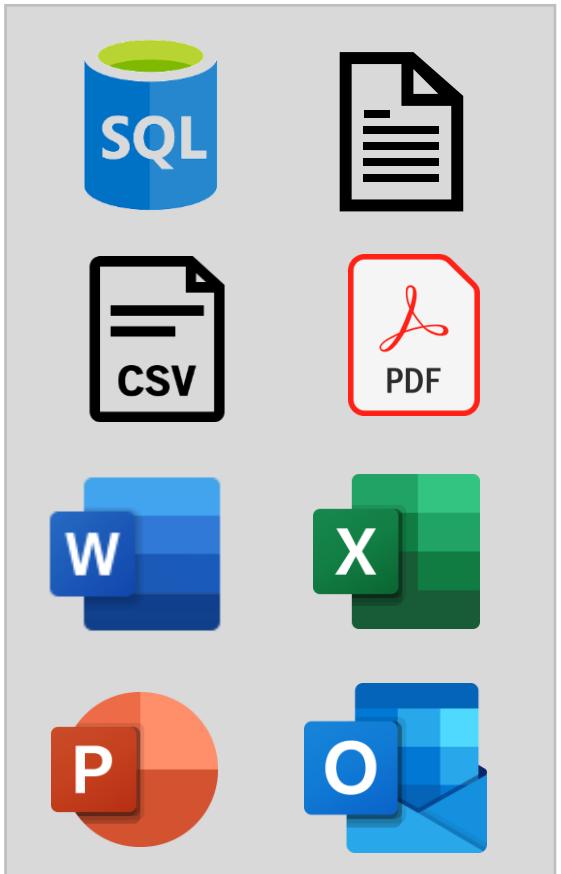


## Search Index

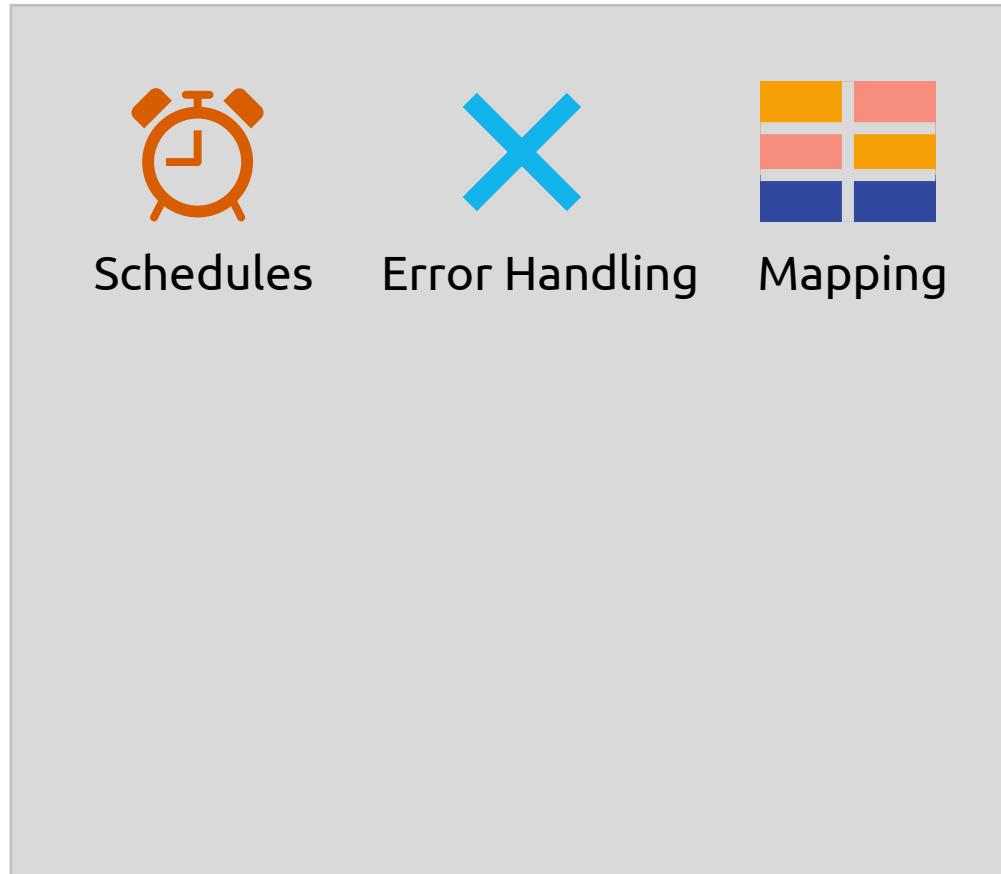
```
{  
  "DocId": 1,  
  "Content": ...  
},  
{  
  "DocId": 2,  
  "Content": ...  
},  
{  
  "DocId": 3,  
  "Content": ...  
}
```

# Cognitive Search Building Blocks – Indexer

## Data Source



## Indexer



## Search Index

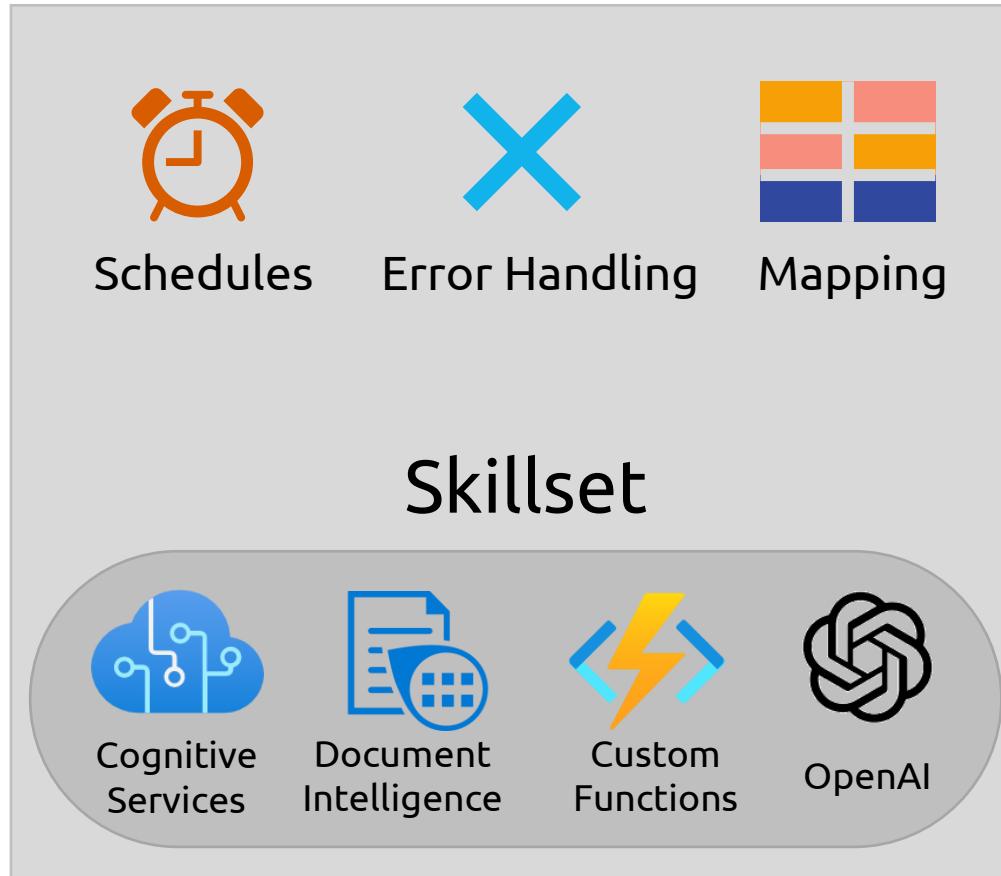
```
{  
  "DocId": 1,  
  "Content": ...  
},  
{  
  "DocId": 2,  
  "Content": ...  
},  
{  
  "DocId": 3,  
  "Content": ...  
}
```

# Cognitive Search Building Blocks – Skillsets & Skills

## Data Source



## Indexer



## Search Index

```
{  
  "DocId": 1,  
  "Content": ...  
},  
{  
  "DocId": 2,  
  "Content": ...  
},  
{  
  "DocId": 3,  
  "Content": ...  
}
```

# Cognitive Search Building Blocks – Skillsets & Skills

{

SkillSet Metadata

List of Skills to apply in Order

Cognitive Services Configuration

Knowledge Store Configuration

}

# Cognitive Search Building Blocks – Skillsets & Skills

{

SkillSet Metadata

List of Skills to apply in Order

Cognitive Services Configuration

Knowledge Store Configuration

}

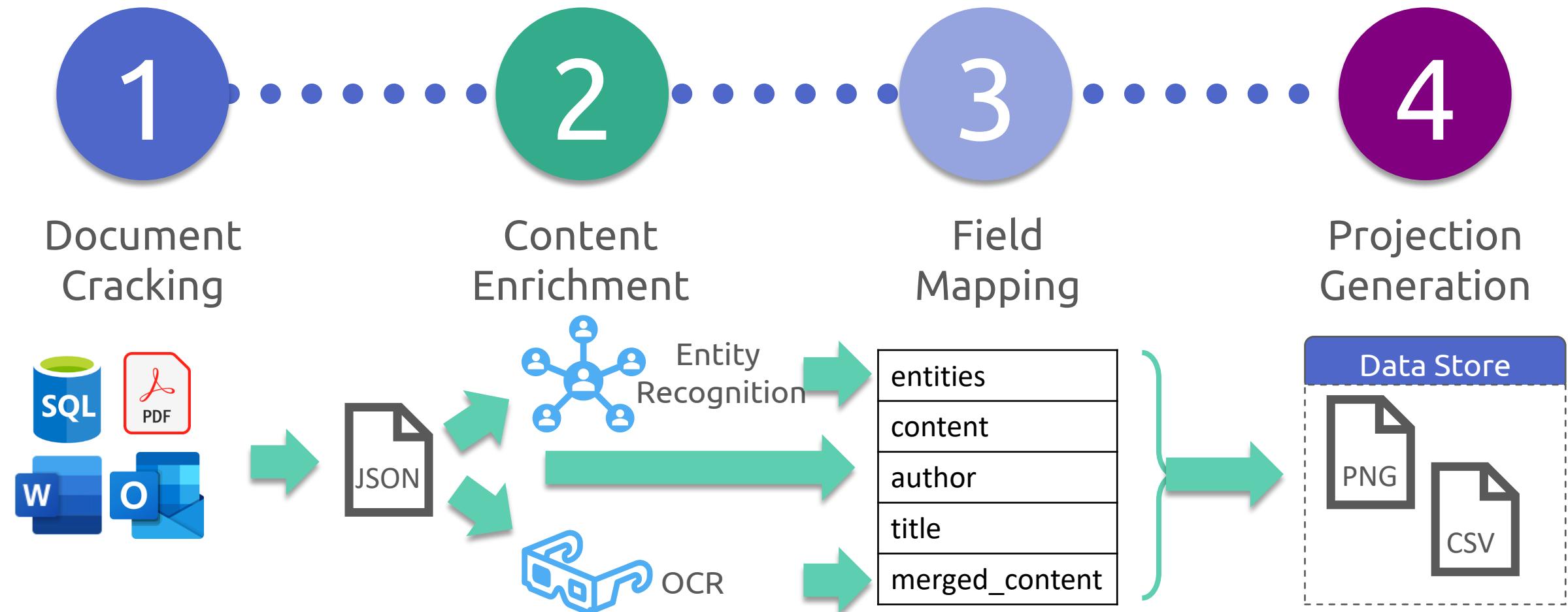
Content Extraction

Optical Character Recognition

Merge OCR and Text

Project Data into a Data Lake

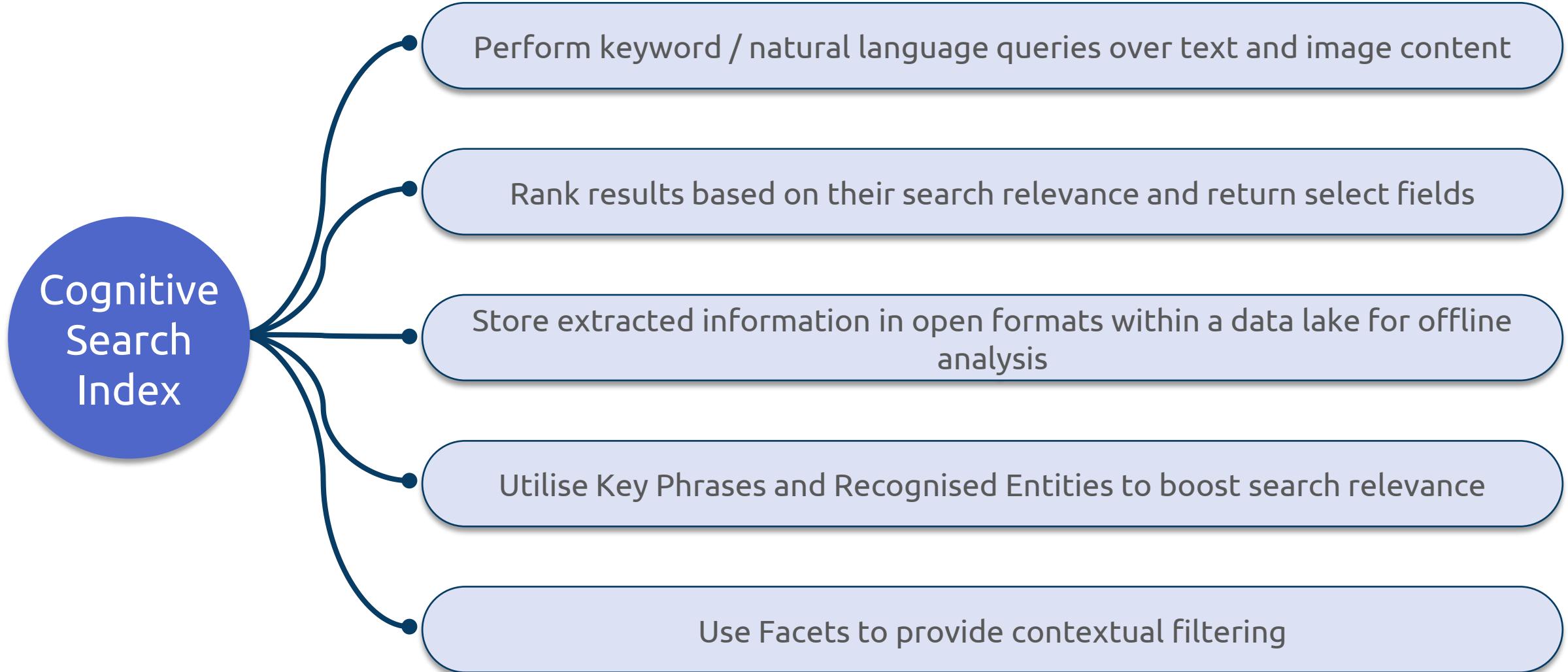
# Cognitive Search Processing Steps



# Cognitive Search Demo:

## Creating and populating an Index

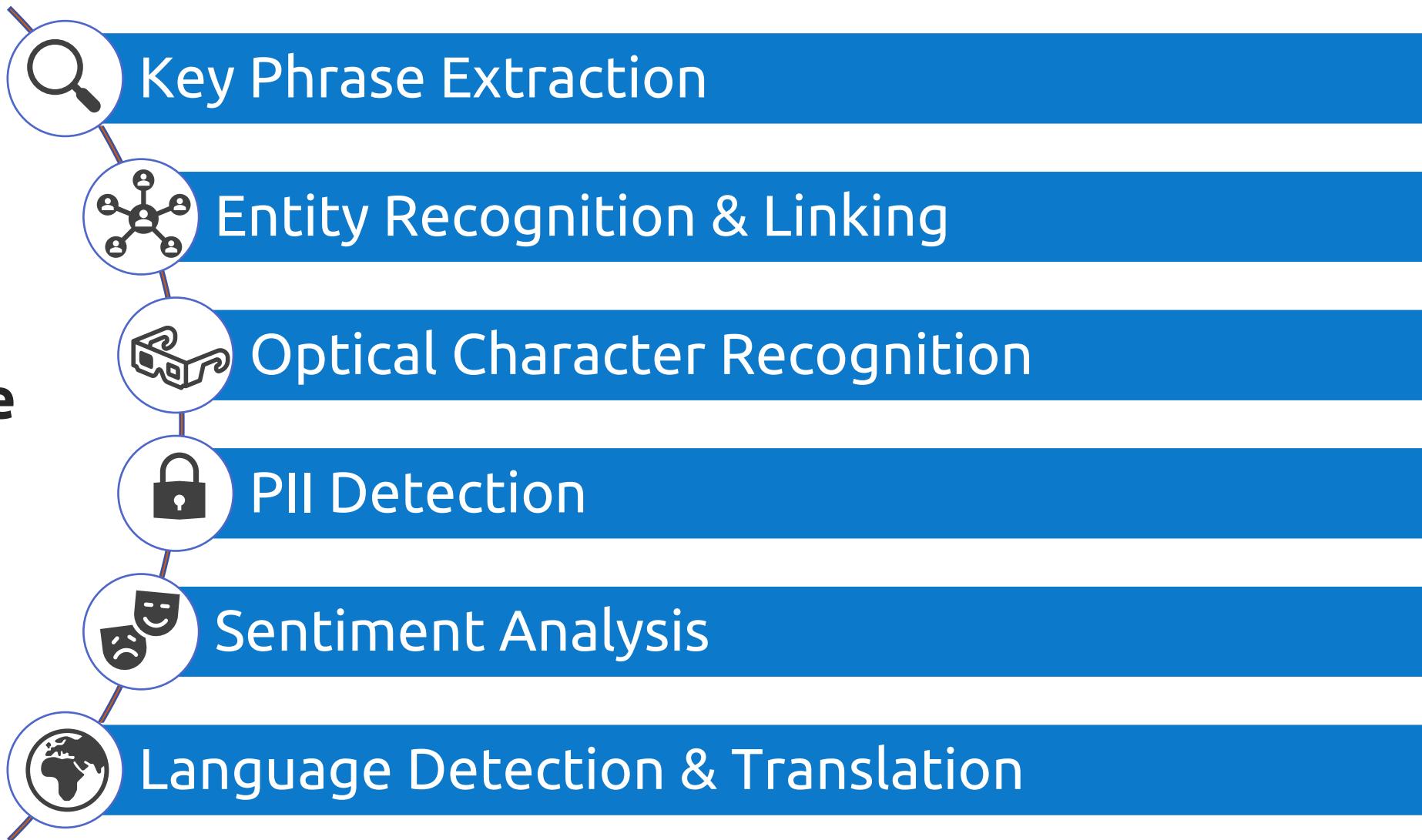
# What can we do now...?



# Custom Skills & AI Enrichments

# Built-In Cognitive Skills

## Built-In Cognitive Skills



# Built-In Cognitive Skills Key Facts

---



Cognitive Skills rely on a connection to an existing Azure Cognitive Service resource



Pricing is based on the specific API usage, generally per 1k transactions or 1k text records.



Each Cognitive Skill will have different outputs depending on the service.

# Built-In Cognitive Skills

```
2 
3     "@odata.type": "#Microsoft.Skills.Text.V3.SentimentSkill",
4     "context": "/document",
5 
6     "includeOpinionMining": true,
7     "defaultLanguageCode": "en",
8 
9     "inputs": [
10        {
11            "name": "text",
12            "source": "/document/content"
13        },
14        {
15            "name": "languageCode",
16            "source": "/document/languageCode"
17        }
18    ],
19 
20     "outputs": [
21        {
22            "name": "sentiment",
23            "targetName": "sentiment"
24        },
25        {
26            "name": "confidenceScores",
27            "targetName": "confidenceScores"
28        },
29        {
30            "name": "sentences",
31            "targetName": "sentences"
32        }
33    ]
34 }
```

**Skill Metadata:** What skill is being invoked

**Parameters:** Any skill specific configurations

**Inputs:** The input values for this specific skill

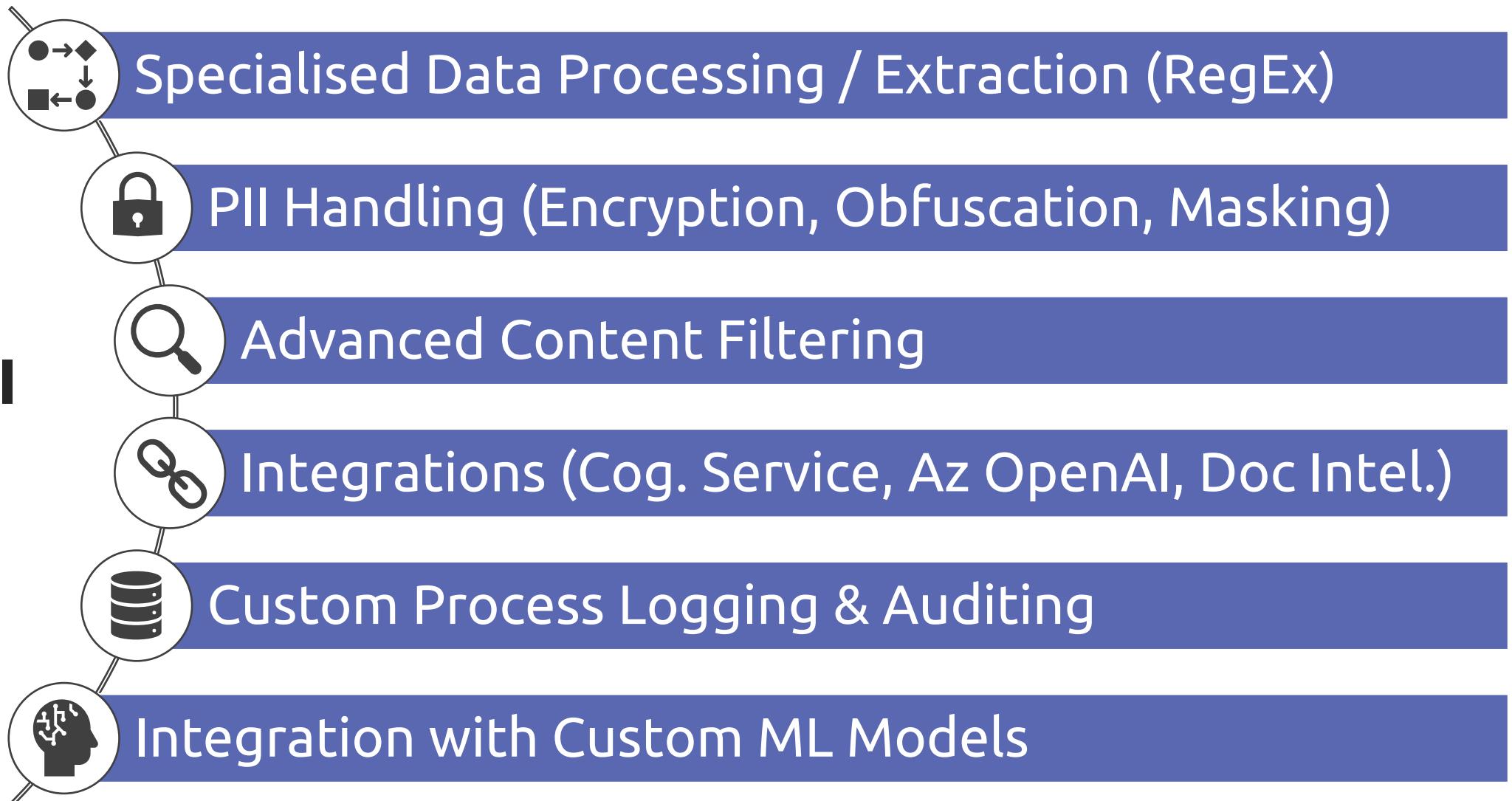
**Outputs:** The expected AI enrichments to be mapped into the Index

# Cognitive Search Demo:

Augmenting documents with Cognitive Services  
(Entity Linking)

# Custom Skill - Web API Skills

## Custom Web API Skills



# Custom Skills – Web API Skills Toolkit



## Azure Functions

To host your API Endpoint in the cloud on a consumption basis



## POSTMAN

For testing your API Endpoint locally and in the cloud deployment



## Visual Studio Code

To develop and test your function code prior to deployment into Azure

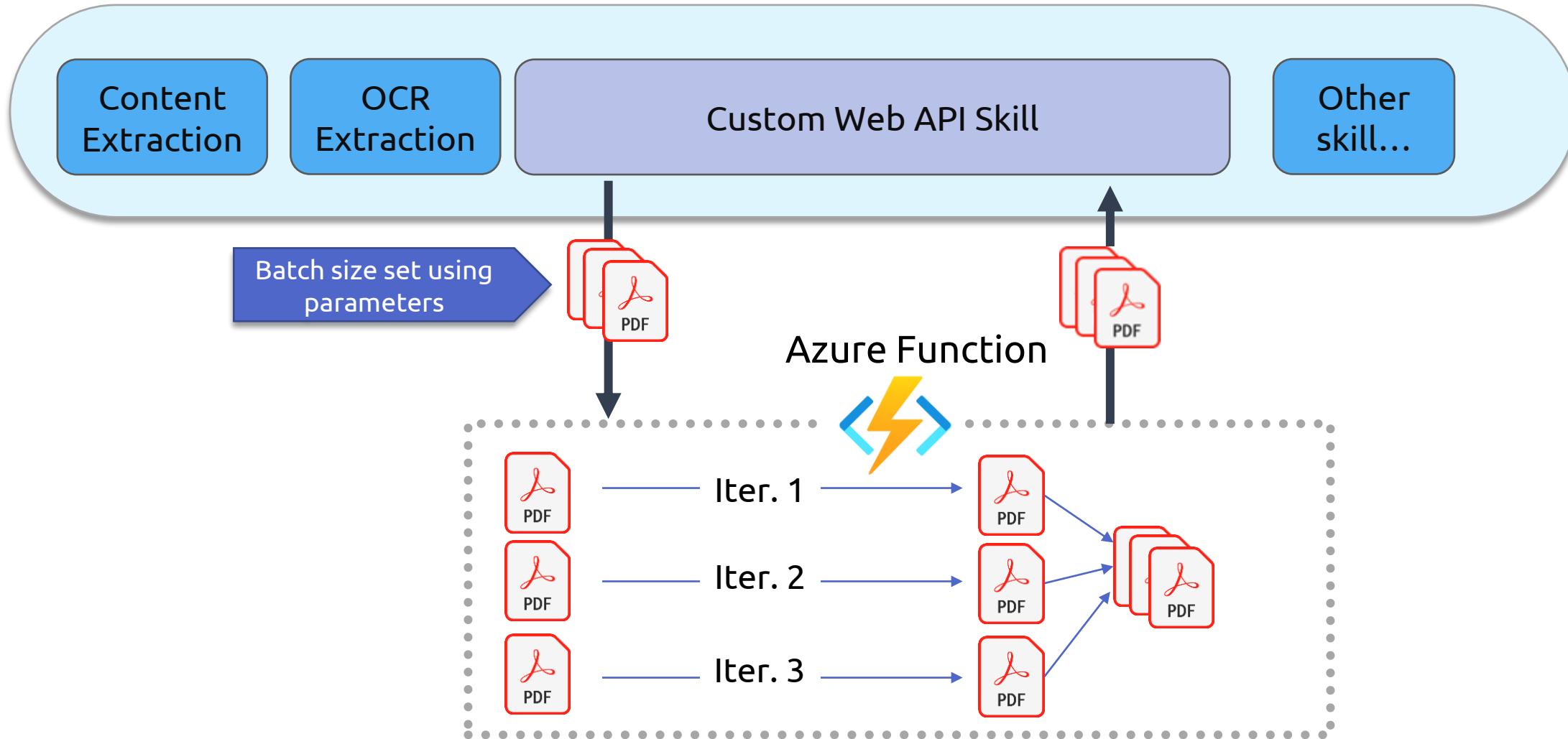


To create the logic to process your document data

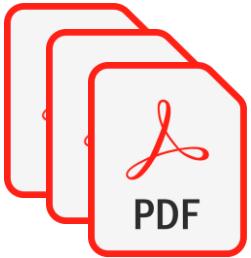


# Custom Skills – Web API Skills

## Cognitive Search Skill Set



# Custom Web API Key Facts



Data will arrive in batches.  
The Custom function needs  
to handle multiple records at  
once i.e using a For Loop

```
{  
  "values": [  
    {  
      "recordId": "0",  
      "data": {  
        "filePath": "path/to/file.pdf"  
      }  
    },  
    {  
      "recordId": "1",  
      "data": {  
        "filePath": "path/to/file.pdf"  
      }  
    }  
  ]  
}
```



Functions can be written in  
C#, Python, JavaScript etc



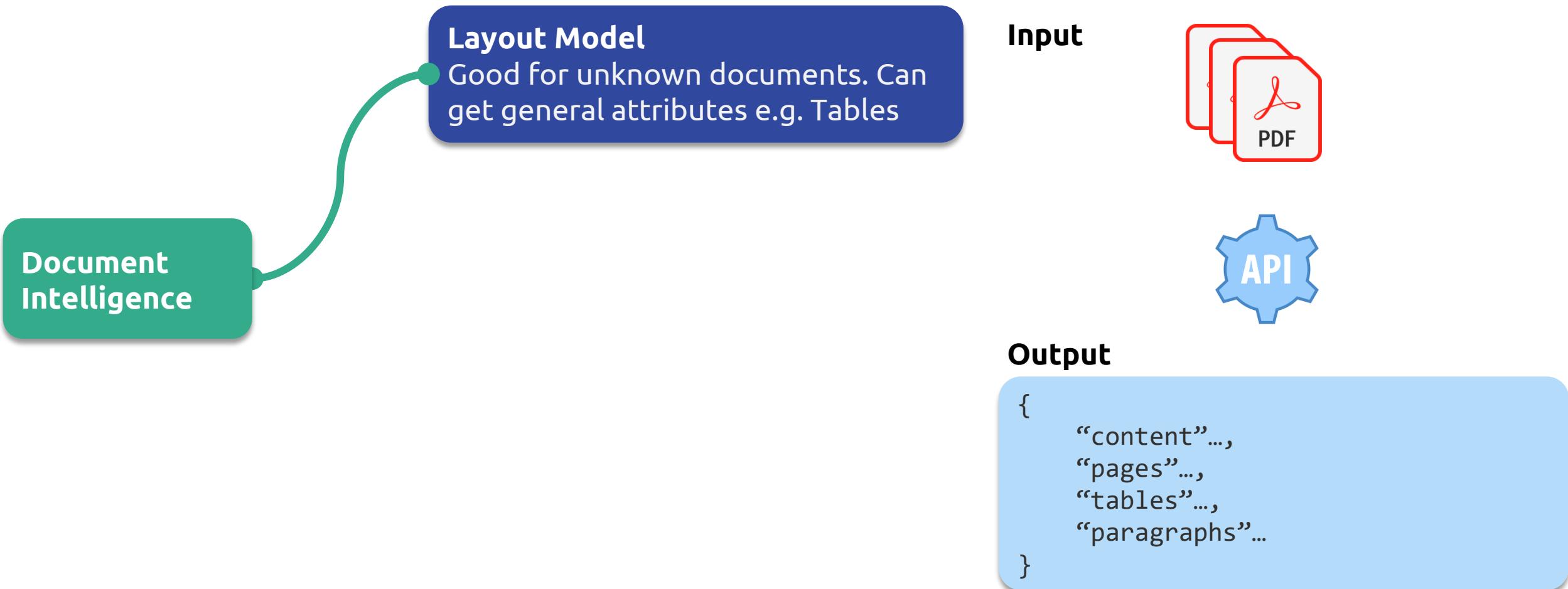
The return JSON must  
contain the same number of  
records in the Batch.  
“try / except” handlers are  
essential!



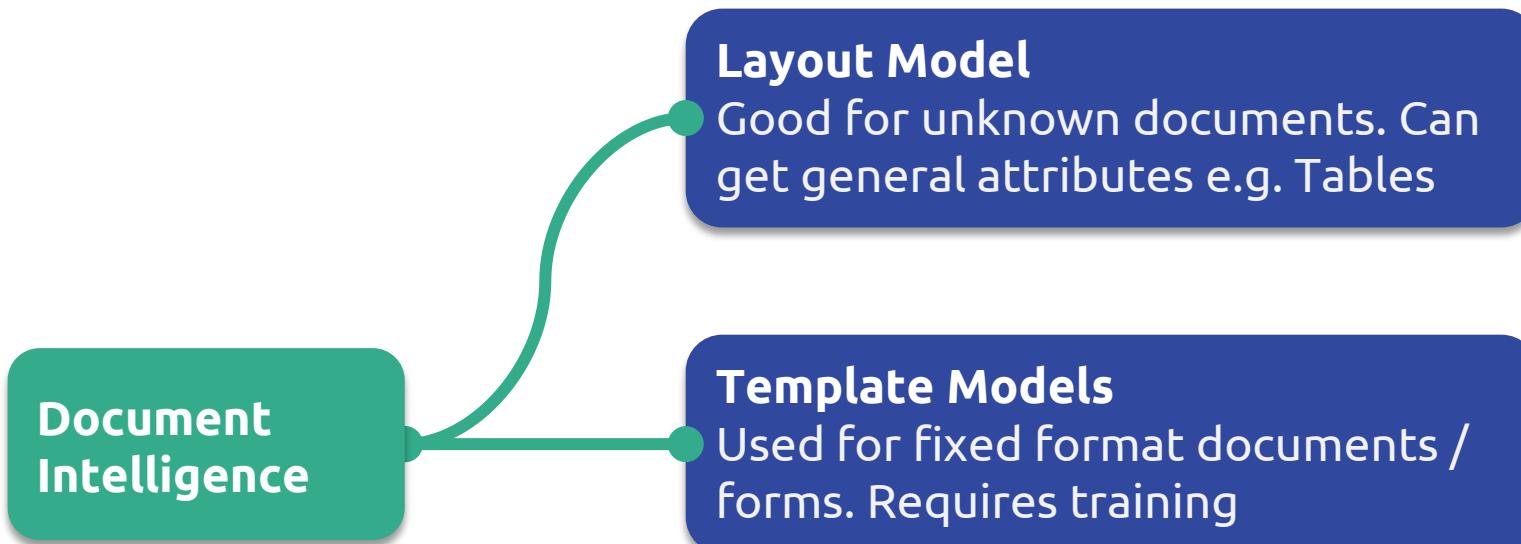
The Custom Function must  
return in the following format:

```
{  
  "values": [  
    {  
      "recordId": "0",  
      "data": {  
        "GPT_Summary": "This is a summary"  
      },  
      "errors": null,  
      "warnings": null  
    },  
    {  
      "recordId": "1",  
      "data": {  
        "GPT_Summary": "Another sample summary"  
      },  
      "errors": null,  
      "warnings": null  
    }  
  ]  
}
```

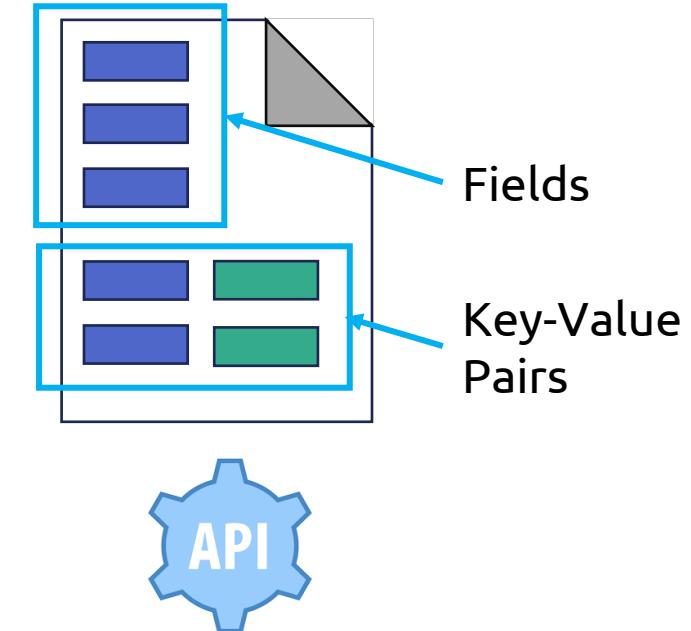
# Azure Document Intelligence (Form Recognizer)



# Azure Document Intelligence (Form Recognizer)



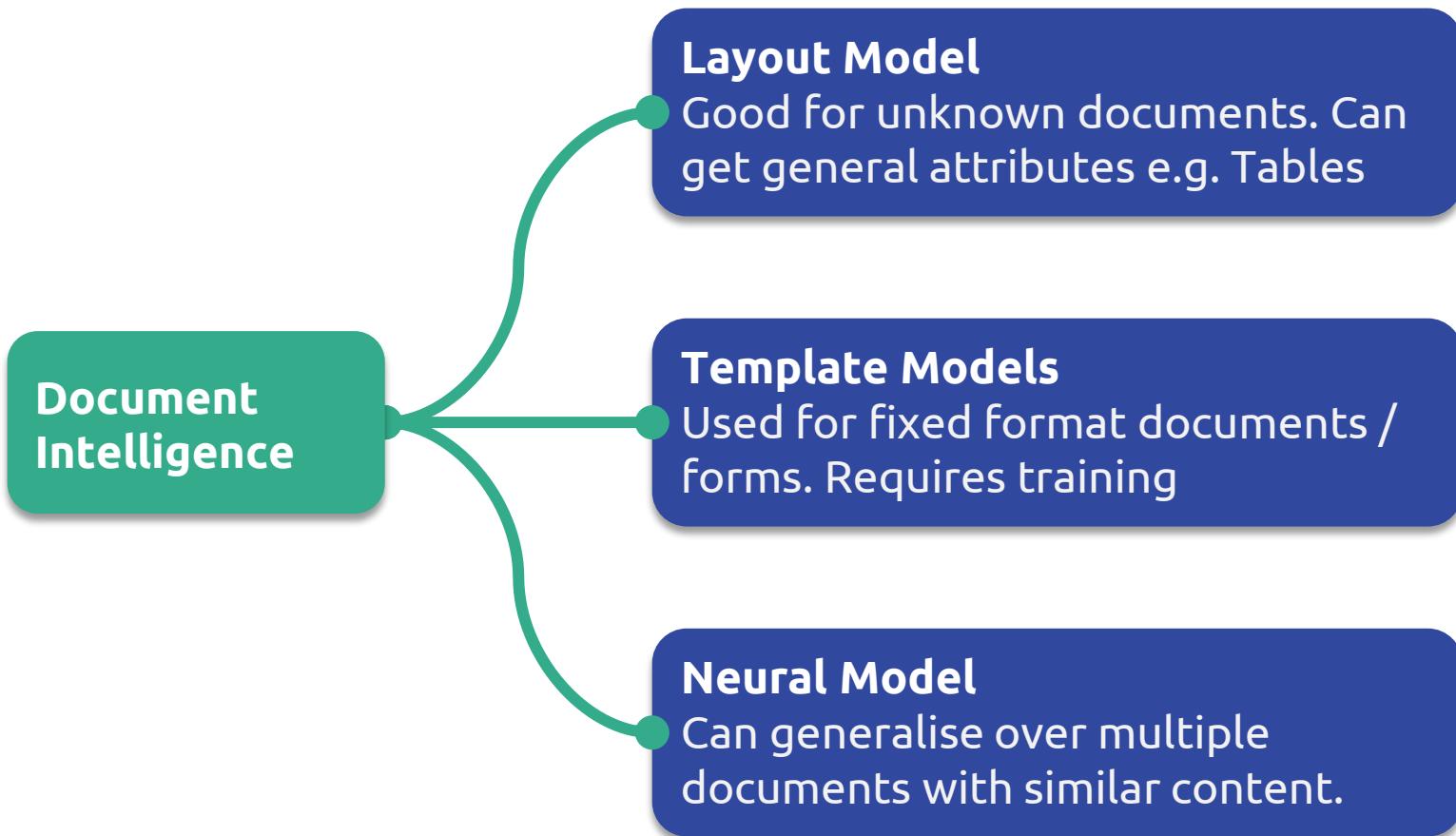
Input



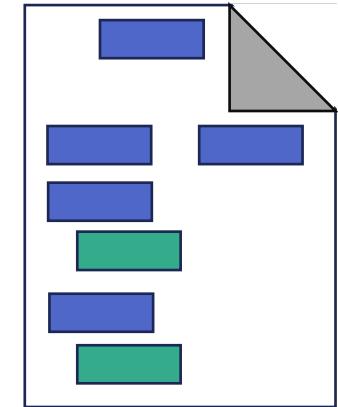
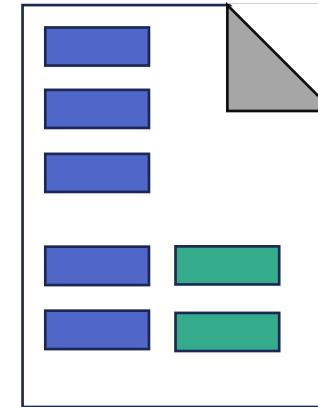
Output

```
{  
  "fields"...,  
}
```

# Azure Document Intelligence (Form Recognizer)



**Input**



**Output**

```
{  
  "fields": ...  
}
```

# Azure OpenAI Endpoints

Azure OpenAI

## Completions API

Generate a predicted response based on the user provided prompt.

## Prompt

Summarise the following text:

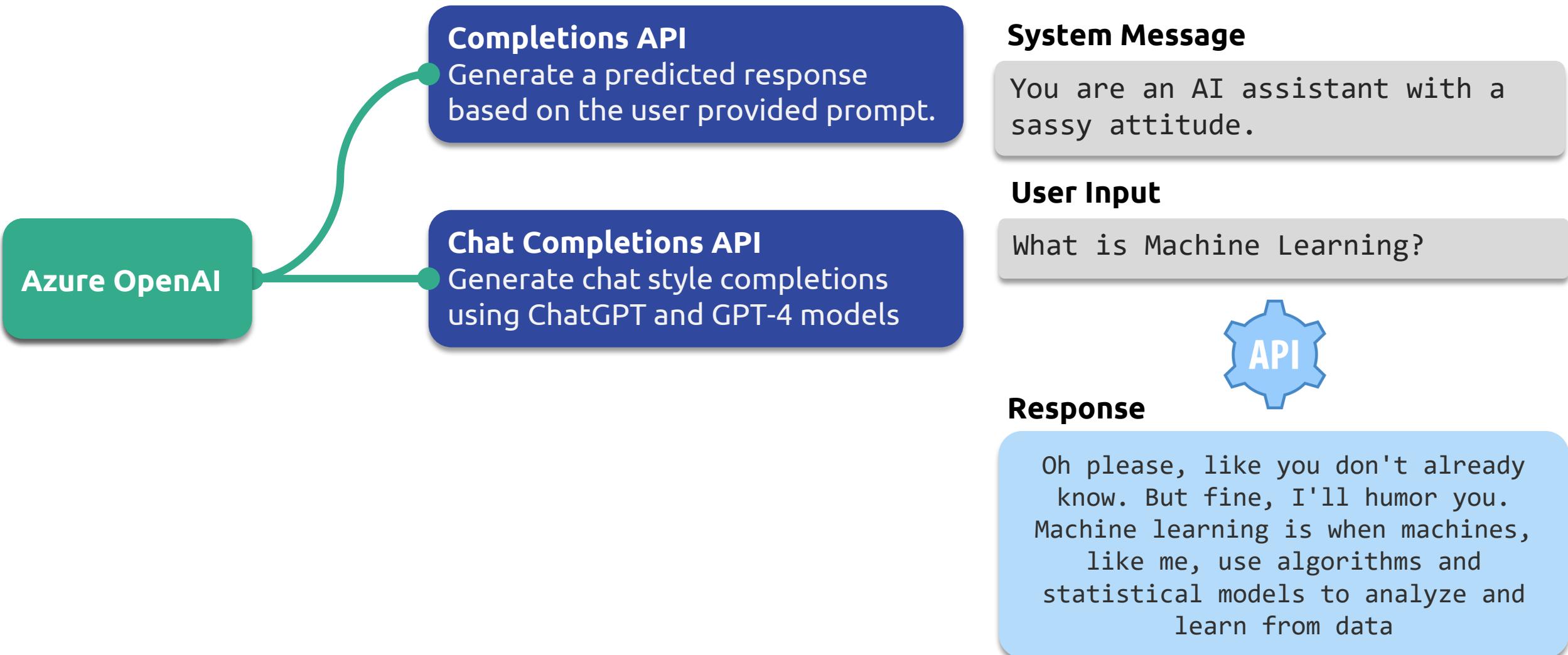
“Machine learning (ML) is a field devoted to understanding and building methods that let...”



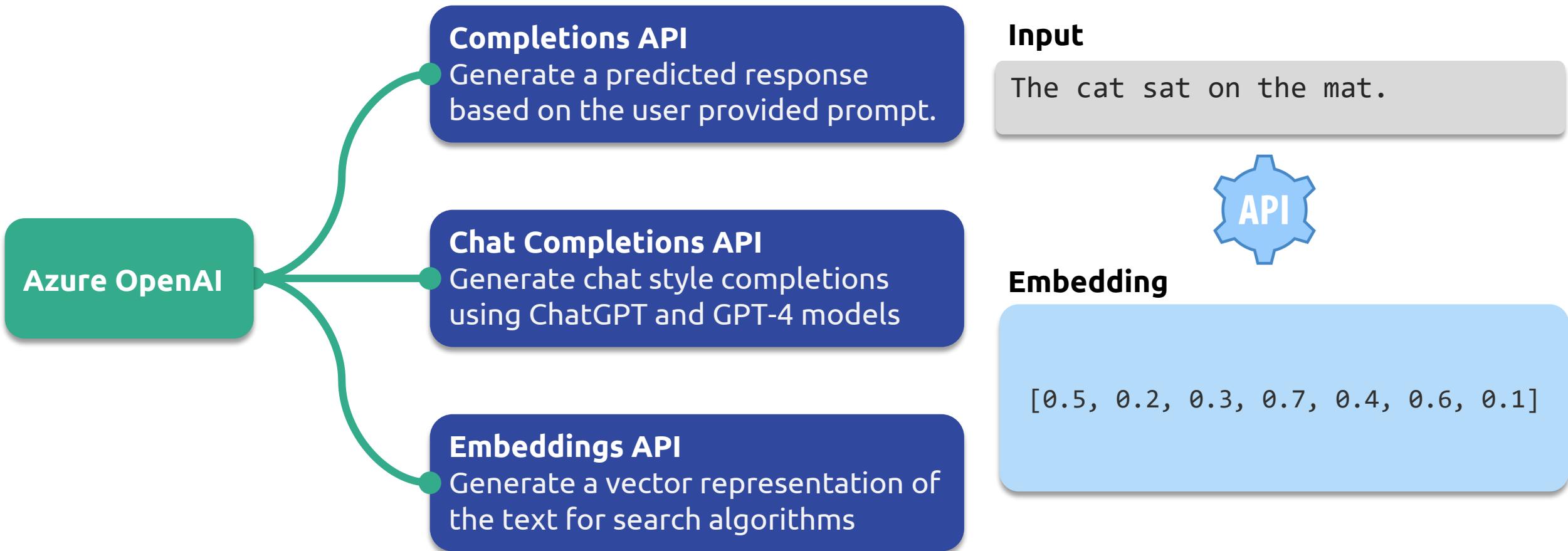
## Response

Machine Learning is about building models to improve computer performance from data.

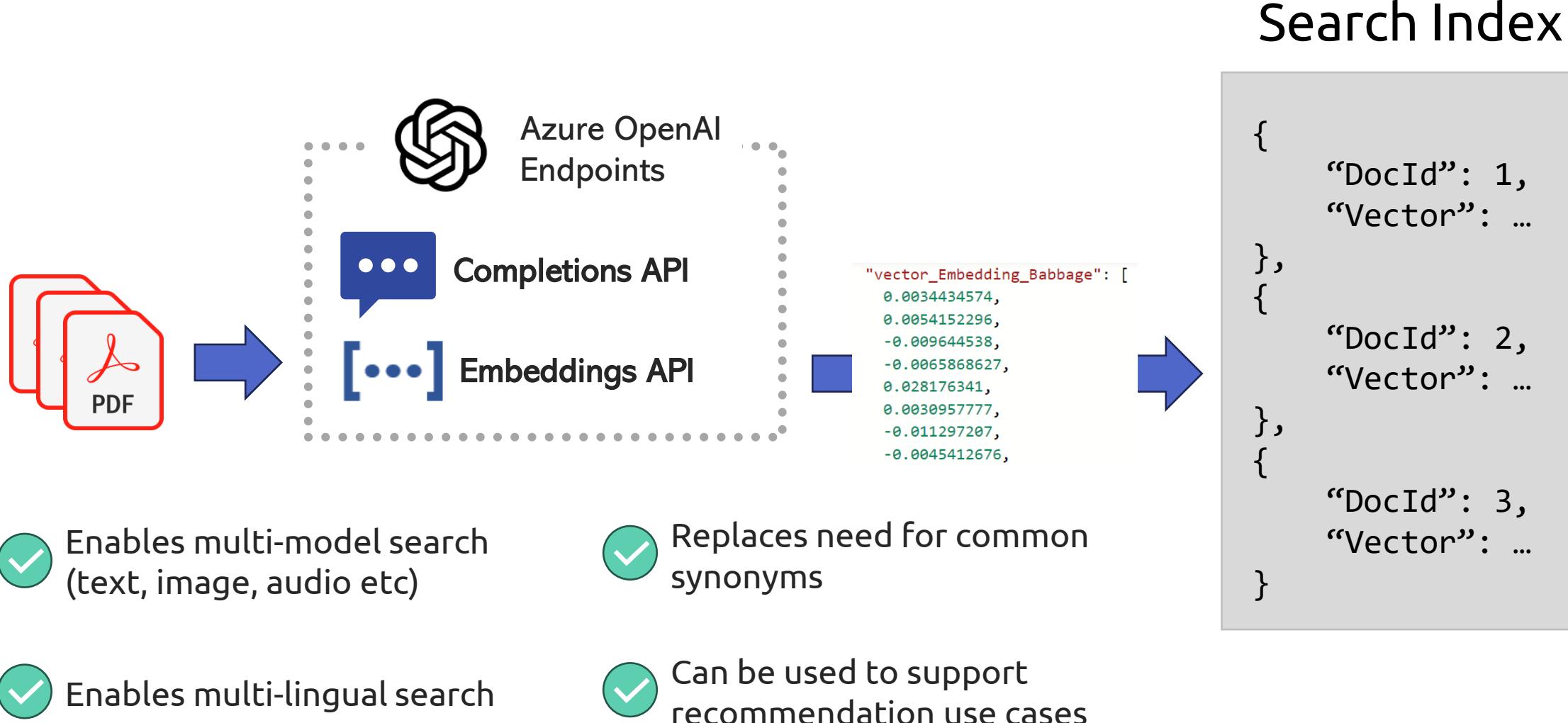
# Azure OpenAI Endpoints



# Azure OpenAI Endpoints



# Custom Skills – Integrate OpenAI & Document Intelligence



# Cognitive Search Demo:

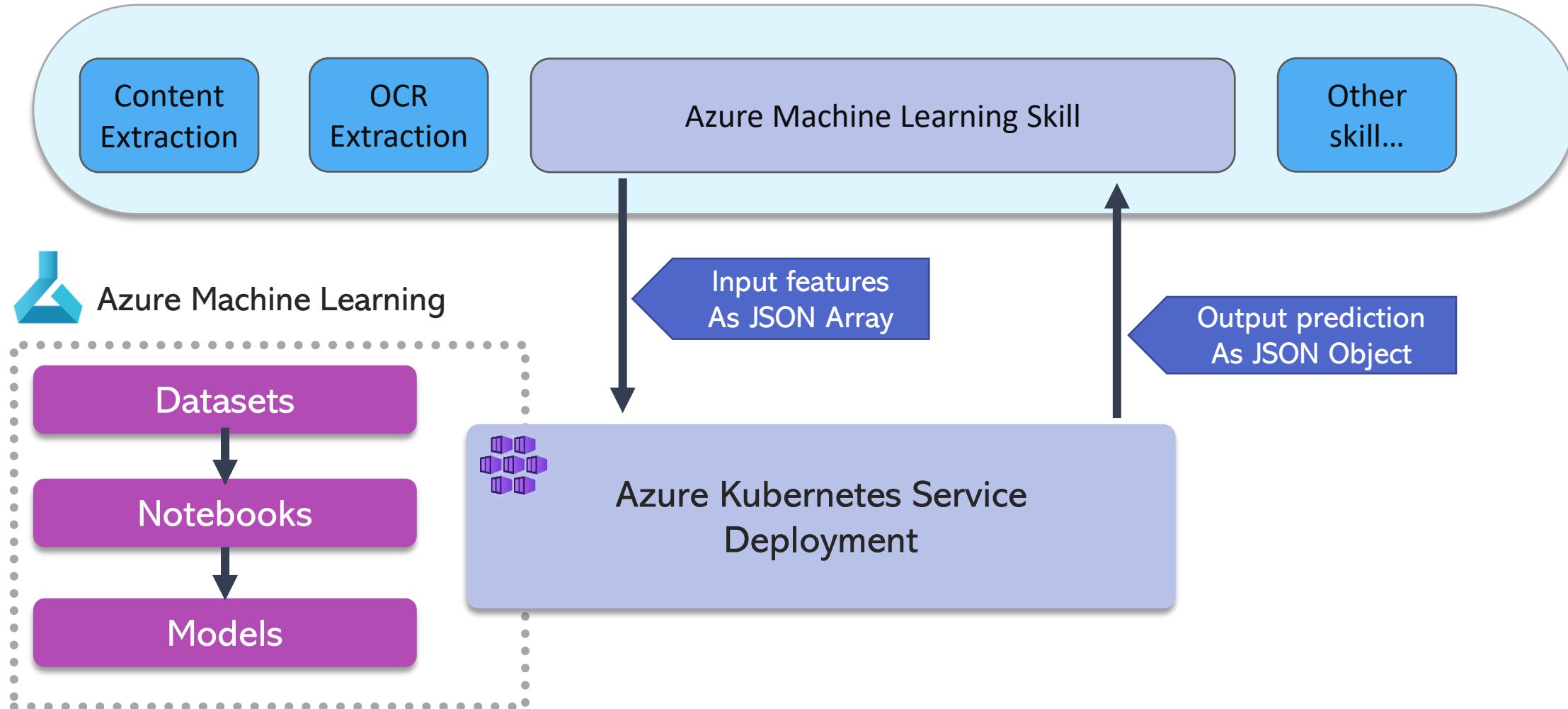
Augmenting documents with Custom Skills  
(Document Intelligence and Azure OpenAI  
Integration

+

Vector Search using OpenAI Embeddings)

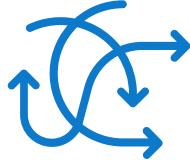
# Custom Skills – Azure Machine Learning

## Cognitive Search Skill Set



# Enhancing Cognitive Search Usage

# Azure Cognitive Search Debug Sessions



Stringing together many AI services gets confusing!



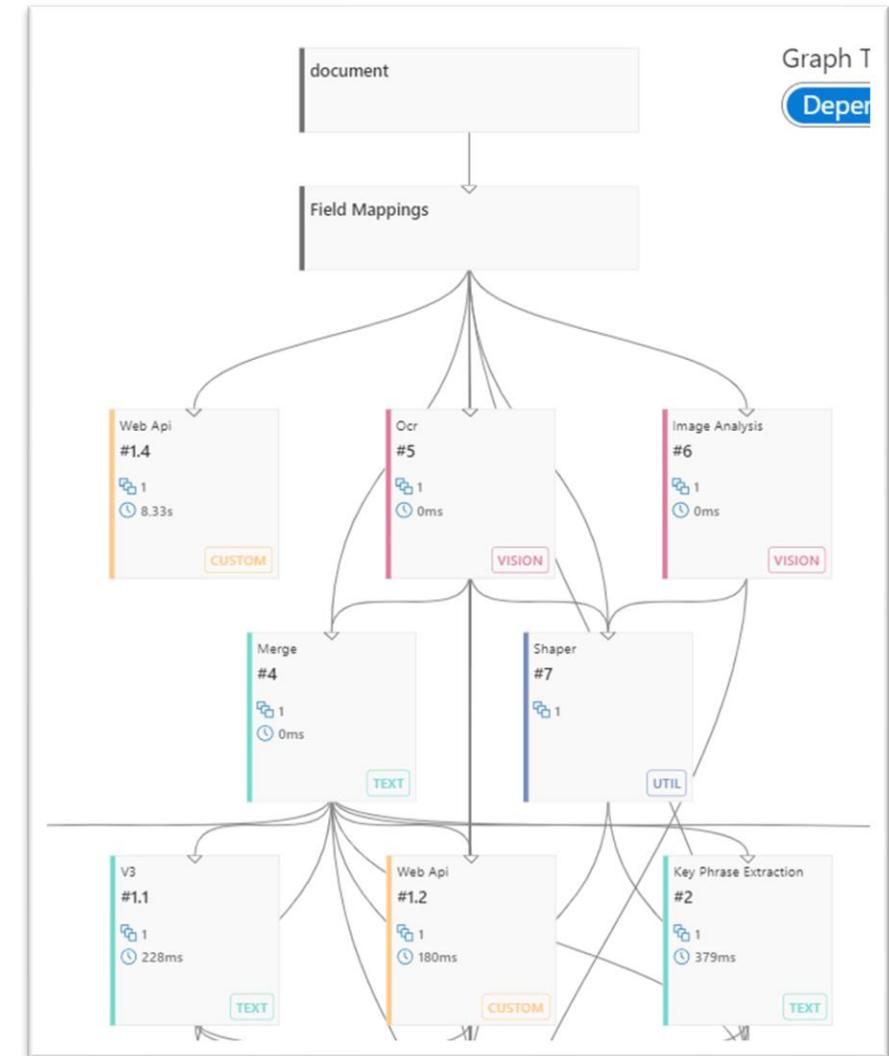
Raw JSON makes it hard to see dependencies



Understanding the enriched structure is tricky!



Cognitive Search has the answer!



# Scoring Profiles & BM25 Tuning



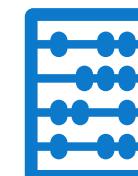
## Scoring Profile

Boost documents that match the search string in each searchable field

Use functions to boost documents based on proximity or magnitude of a field

BM25 parameters can be configured on the index based on testing

k1 defaults to 1.2. Can be between 0 & 3  
B defaults to 0.75. Can be between 0 & 1



## BM25 Tuning

k1

Allows the search to prioritise documents that match more terms from the search

b

Controls how much the length of a document effects the final score. E.g. Allows us to penalise longer docs

# Semantic Search & Semantic Configurations



Allows more powerful searches to be conducted, based on Semantic models – not just lexical



Can be enabled via Portal and toggled on and off for each query



Semantic Answer  
Generation  
(Extractive)



Semantic re-  
ranking over  
search score



Configured using  
a Semantic  
Configuration



Operates over the  
first 50 results  
from BM25 search

# Cognitive Search Demo:

Cognitive Search Enhancers

(Debug Session

+

Scoring Profile

+

Semantic Config)

# Additional Tools

<https://github.com/Azure-Samples/azure-search-power-skills>

A screenshot of a GitHub pull request page. The pull request is titled "Merge pull request #129 from Azure-Samples/dependabot/pip/Template/Py...". It shows 11 commits across various branches and files. The commits are:

- .github Initial commit
- Common Bump Newtonsoft.Json from 12.0.3 to 13.0.1 in /Common
- Geo/GeoPointFromName Complete the transition to 'main' name for the default branch
- SampleData Add analyze form skill
- Template Bump mlflow from 2.5.0 to 2.6.0 in /Template/PythonFastAPI
- Tests/PowerSkillTests Reverted Missing endpoint to being an error
- Text Bump transformers in /Text/TextSummarization/powerskill
- Utils Complete the transition to 'main' name for the default branch
- Vector/EmbeddingGenerator Merge pull request #125 from hyoshioka0128/main-1
- Video/VideoIndexer Complete the transition to 'main' name for the default branch
- Vision Bump requests in /Vision/AutoMLVisionClassifier/notebooks

<https://github.com/Azure-Samples/azure-search-lab>

A screenshot of the Azure Search Query Composer interface. The interface is a web-based tool for building search queries. It has several sections:

- Parameters:** Request Method: GET, API Version: 2023-10-01, Search Fields: content, Search Text: \*, Search Mode: any, Select Fields: content, Query Type: simple, Session id: session id, Minimum Coverage: 100, Skip: 0, Top: 100, Count: false.
- Result:** Order by: metadata\_storage\_size asc, Search Score asc, Geo distance Point X Point Y asc.
- Highlight:** content, Highlight PreTag: </em>, Highlight PostTag: </em>.
- Filters:** filters expression, Learn more about filters.
- Scoring Statistics:** local, Scoring Profile: , Scoring Parameter: scoringParameter.
- Facet:** metadata\_content\_type count count number.
- Semantic Configuration:** Without-Semantic-Config, Query Language: English.
- Speller:** lexicon, Answers: extractive, Captions: extractive.

The bottom section shows the generated query string: &api-version=2023-10-01&search=&searchMode=any&queryType=simple&\$skip=0&\$top=100&\$count=false&scoringStatistics=local&minimumCoverage=100&\$orderby=search.score() asc. A "Query" button is at the bottom right.

# Questions & Feedback

[www.adatis.co.uk](http://www.adatis.co.uk)

[matt.how@adatis.co.uk](mailto:matt.how@adatis.co.uk)



