# Cisco Document Intelligence Final Presentation

**MLDS Team: Andrew, Kelvin, Xin, Xingyu, Ziqiao**

Northwestern

# Agenda

# Problem Definition

**Problem Statement:**

- Cisco supplier product specification files are unstructured
- Product attributes are not consistently formatted or organized within documents
- Current data extraction process is extremely manual and time consuming

**Objective:**

- Establish an automated data extraction pipeline to collect the information needed from unstructured document sources
- Achieve high accuracy and adaptability to different types of files in the extraction process
- Build an user-friendly UI that non-technical people can interact with
- NOTE: this project will focus on key attribute extraction text-rich PDF files obtained from different suppliers

# Example PDF Documents

Document link: *RK73H.pdf*

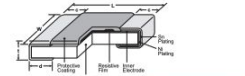Document link: Gs-12-1565.pdf

Document link: *Vzh680m2atr-1313s.pdf*

# Data Overview

1.  **Attributes have been tested (10 in total):**

    *Extracted attributes will depend on product type. Testing was done with the attributes below.*

    a.  Supplier Name, Product Type, Dimensions, Orientation if any, Current Rating, Voltage, Frequency, Impedance, Capacitance, Temperatures

2.  **Selected sample PDF documents (9 in total):**

    *All selected products are commodity type for confidentiality reasons and for experimentation purposes*

    a.  **Capacitors:**

        i.  ***KGM15CR51E106K-DATA.pdf***   ***Vzh680m2atr-1313s.pdf***

    b.  **Molex:**

        i.  2031430001-PS.pdf

    c.  **Amphenol:**

        i.  Ssio_mini_cooledge_0_60mm.pdf      Gs-12-1457.pdf      Gs-12-1565.pdf

    d.  **Resistance:**

        i.  PE_Lseries.pdf         ***RK73H.pdf***

    e.  **App Processor:**

        *i.*  ***IMX8MMIEC.pdf***

**Note:** we selected 9 PDF documents with a diverse range of vendors, products, and document structures. PDFs highlighted in bold and italic are the ones with ground truth values and thus are mainly used for evaluation purposes.

# Solution Methodology Overview

- How would we have solved this problem traditionally?
  - A **rule-based/pattern matching approach** (such as RegEx); however, this would have had to be tailored to a specific manufacturer and document type
  - **Machine Learning** based methods (such as SVM, CNN or LSTM) could have provided some generalizability; however, it would require a massive exercise in generating a training data with accurate ground truths

- Why is Generative AI recommended for this use-case?
  - Recent breakthroughs in NLP, particularly with transformer-based models like GPT-4, have significantly **improved the understanding and generation of human-like text**
  - These models have an **enhanced ability to understand context**, which is crucial in accurately extracting relevant information from documents.

# Solution Methodology Overview (cont.)

- **What are the advantages/disadvantages of GenAI ?**
  - **Significant Advantages**:
    - *Adaptability and Generalization*: GenAI models are not as rigid as rule-based systems and can generalize better than traditional machine learning models
    - *Contextual Understanding*: GenAI models are better at understanding context and can understand nuances in language and structure that rule-based systems might miss
    - *Multimodal*: new GenAI models are multimodal meaning that they can take in both text and images
  - **Significant Disadvantages**:
    - *Complexity and Resource Intensity*: GenAI models are complex and often require substantial computational resources.
    - *Interpretability and Explainability*: GenAI models are often described as "black boxes" due to their lack of interpretability. Understanding why a model made a particular decision can be challenging, which is problematic in scenarios where transparency is crucial.

# Solution Workflow



Product Spec Documents → PDF Extraction Module → Text-based document → LLM (w/ RAG) → Meaningful Product Attributes

**Meaningful Product Attributes**

Attributes can be stored in a database, displayed on a website, summarized in a table, etc.

Takes a photo of the document and then uses optical character recognition to turn the photo into text

Interface with LLMs via "prompts" or requests for action

# PDF Extraction Module: PDF Pre-Processing

## PDF Extraction Methods Overview

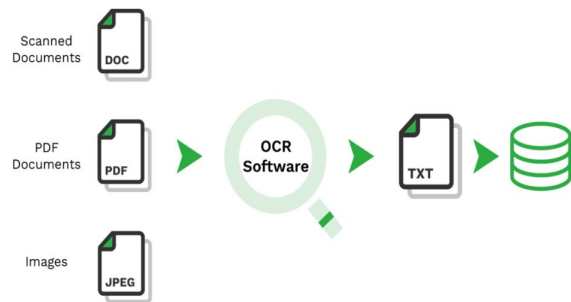PDF extraction is the process of retrieving text, images, data, and metadata from PDF documents using specialized software or techniques, enabling the conversion of content from a non-editable format into a usable, editable format for various applications.

## OCR (Optical Character Recognition)-Tesseract
Extracts printed or written text from images. It was originally developed by Hewlett-Packard, and development was later taken over by Google.



Use **Pdf2Image** first convert pdf files to images, then use OCR to extract texts from images.

*Input pdf or image*



*"Image word"*

*Output a single column of texts*

*able to extract*

```
1   KYOCERA
2
3   Dimension
4   1608(JIS) / 0603(EIA) L ~ w ;
5   Unitmm _ a DC Bias AC Voltage
6   20 20
7   0 0
8   Nut −20 _ −20
9   x. = 40
10  Characteristic O 40 oO
11  Temp. Range −55deg to85deg GS −60 O −60
12  Temp. coeff +/−15% −80 −80
13  Capacitance 10uF −100 es errors −100
14  Tolerance +/−10% 0 5 10 15 20 25 0.001 0.01 0.1 1
15  Rated Voltage 25Vdc DC BIAS (V) AC voltage (Vrms)
16  S parameters(Series) Impedance/ESR Temperature Characteristics
17  −_−S11 am S71 −=−Z(0Vdc) ==Rs(0Vdc) am−OVdo me 12.5Vde
18  10 20
19  0
20  0
21  1
22  −20
23  ao E −20
24  = −40 s iS
25  a Sf 0.1 XS −40
26  g o o
27  ~ 60 Wu S −60
28  " N 0.01
29  −80 −80
30  −100 0.001 −100
31  0.01 0.1 1 10 100 1000 0.01 0.1 1 10 100 1000 −60 −40 −20 0 20 40 60 80 100
32  Frequency (MHz) deg
33
34  Frequency (MHz)
35
36  Capacitor Division
```

# PDF Pre-Processing: PDF Extraction Methods

| Technique/Model | Pros | Cons |
|---|---|---|
| **pdf2image + OCR-Tesseract** | • Open-source and free (can be used locally)<br>• Good accuracy and versatility<br>• Can read "image words" | • Requires image preprocessing to achieve optimal results<br>• Outputs a single column of text, does not preserve the structure |
| **Pdfminer/Pdf Plumber** | • Simple and robust; preserves control characters such as "\n" | • Cannot read the "image words" such as vendor name |
| **PyPDF (PyPDF2& PDFLoader)** | • PDFLoader has more compatibility with LangChain and so does PyPDF2 with LLaMa | • Some signs cannot be read correctly. E.g. Read °C as oC<br>• Cannot read "image words" |
| **GPT4.0-image reading** | • Good accuracy and versatility<br>• Can read "image words" | • No accessible API, not free |

# Large Language Models (LLM)

**What is a Large Language Model?**

- A large language model (LLM) is a type of artificial intelligence that learned on massively large data sets to understand, summarize, generate and predict new content.

**How to interact with a Large Language Model?**

- Ask clear and specific questions with possible context information
- Follow up on previous questions to improve the answer or ask for new related questions
- **Example use case:** use prompts for question answering, creative tasks like writing stories, poems, or even generating ideas, and more

# Large Language Models (LLM)

**Open-source vs Closed-source**

- **Open-Source Models:** publicly accessible and can be modified or distributed by anyone,   can be used without restrictions. Examples: Mistral-7B, Mosaic-mpt7b.
    a.   **Pros:** Free to use, customizable, community-driven improvements.
    b.   **Cons:** Less support, can require significant resources to train and deploy, potentially lower performance compared to state-of-the-art closed-source models.
- **Closed-Source Models:** These are proprietary models developed by organizations and are not publicly accessible for modification. Example: GPT-3.5, GPT-4.0.
    a.   **Pros:** Often more advanced, better performance, reliable support and updates, easier to implement.
    b.   **Cons:** Costs (licensing or usage fees), less customizable, reliance on a third-party provider.

**Cost Estimates**

- **Open-Source Models:** Mostly free, but include potential costs for training (compute resources), maintenance, and integration.
- **Closed-Source Models:** Typically involve licensing fees or pay-per-use models. Provide general cost estimates (if available) or ways to access cost information (e.g., contacting providers for quotes).

# Prompt Engineering - The Basics

Prompt engineering is the practice of **carefully crafting and optimizing input text** to effectively guide an AI model to produce desired outputs. Benefits of Prompt Engineering include:

- **Enhanced Accuracy**:
    - Modify the prompt to include detailed explanations for each attribute
    - Guide the LLM to provide comprehensive answers that cover all relevant aspects of an attribute but avoid redundancy
    - Rather than asking about all attributes in a single prompt, call the LLM individually for each attribute
- **Avoid Hallucination**:
    - Include a system message in the prompt instructing the LLM to respond with "I don't know" when it is unsure
- **Correct format**:
    - Provide instruction to the LLM about JSON format and give it an output example. Write a nested dictionary if there are multiple values under one attribute

**In summary**, revise prompts for improved LLM performance, while considering the effort involved in prompt engineering and factoring in the token costs associated with longer prompts.

# Example Attribute Extraction

**Example 2: Voltage Extraction**

**Voltage:**
**500 volts DC;** 6.3 Dielectric Withstanding Voltage,
**500 volts DC;** 6.2 Insulation Resistance

# Example Attribute Extraction

9.2 Inspection Conditions

Unless otherwise specified herein, all inspections shall be performed under the following an conditions:

a. Temperature: 25 +/- 5 deg C

b. Relative Humidity: 30% to 60%

c. Barometric Pressure: Local ambient

8.3 High Temperature Life –EIA 364-17

a. Method A

b. Test Temperature - 85 deg C

c. Test Duration - 500 hours

8.1 Thermal Shock – EIA 364-32.

a. Number of Cycles - 5

b. Temperature Range : Between -55 and +85 deg C

c. Time at Each Temperature - 30 minutes minimum

d. Transfer Time - 30 seconds, maximum

**Example 1: Temperature Extraction**

25+/-5 degC; 9.2 Inspection Conditions
-55 and +85 deg C; 8.1 Thermal Shock
85 deg C, 8.3 High Temperature Life

Northwestern

# Prompt Engineering - Basic Prompt

| Sample Prompts | Output Example | Observations |
|---|---|---|
| "I'm looking for information in this document.<br><br>Here is the information I'm looking for Supplier Name, Product Name, Product Type, Dimensions, Orientation, Current Rating, Voltage, Frequency, Impedance / Capacitance, Temperatures" | By GPT-4:<br>● The supplier name is Amphenol ICC. The product names are EXAMAX2™ and EXAMEZZ2™ Connector System. The product type is Mezzanine connectors.<br>● The document does not provide specific information on dimensions and orientation.<br>● The operating current rating can be found in the AICC Power Design Guide GS-20-0456.<br>● The agency voltage rating is less than 30 VAC RMS / DC and the non-agency voltage rating is 150 Vac RMS or Vpc Maximum.<br>● The document does not provide information on frequency and impedance/capacitance. The operating temperature range is -55 to 85 °C. | ● Output **texts** with a list of attributes and results.<br>● A slash in the attribute list was not conducive to the final output structure. Sometimes it would separate the two attributes sometimes keep them together. |

# Prompt Engineering - Final Prompt

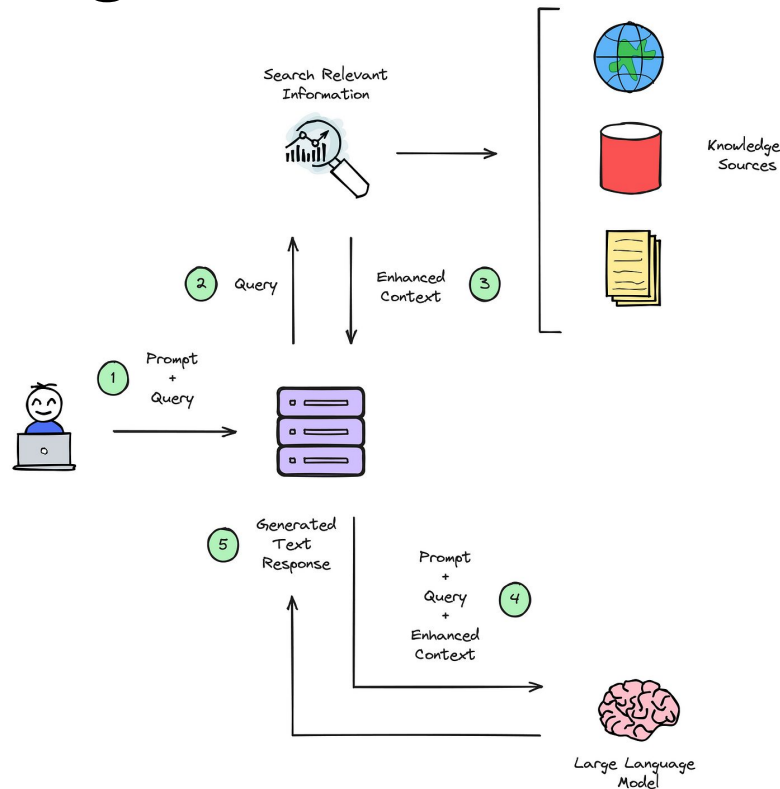| Sample Prompts | Output Example | Observations |
|---|---|---|
| **System message:** Use the following pieces of context to answer the question at the end. If you don't know the answer, just say that you don't know, don't try to make up an answer.<br>**Prompt**: Please extract and compile **all information** related to the following attributes from the document, **avoiding any duplication**. For each attribute, include the specific source or context for each unique value. **Present the findings of each attribute in a single dictionary entry. Exclude any redundant information.**<br>    (a) Supplier Name<br>    ...<br>    (i) Temperatures(temperatures under any conditions)<br><br>The answer should be in the following format:<br>    {<br>      "attribute": "Supplier Name",<br>      "values": [<br>        {"value": "unique value 1", "source": "context or source"},<br>        {"value": "unique value 2", "source": "context or source"}<br>      ],<br>      "summary": "Provide a general summary or observation about the extracted information, focusing on unique aspects."<br>    },<br>    {<br>      "attribute": "Product Type",<br>      ...<br>    },for dimensions, orientation and so on. | By GPT-4:<br>{<br>"attribute": "Temperatures",<br>"values":<br>[{"value": "25 +/-5 degC", "source": "9.2 Inspection Conditions"},<br>{"value": "-55 and +85 deg C", "source": "8.1 Thermal Shock"},<br>{"value": "85 deg C", "source": "8.3 High Temperature Life"} ],<br>"summary": "The document provides information on temperatures under different conditions."<br>},<br>{<br>"attribute": "Voltage",<br>  "values":<br>  [{"value": "500 volts DC", "source": "6.3 Dielectric Withstanding Voltage"},<br>   {"value": "500 volts DC", "source": "6.2 Insulation Resistance"}],<br>"summary": "..."<br>} | ● Output **JSON** with **nested dictionaries** of attribute names, multiple attribute values and source.<br><br>● Both GPT 4 and GPT 3.5 are able to output in the instructed format, and provide reference (e.g. explanation, section number or section title) |

# Multi-Prompt vs Single-Prompt

| Prompt | Definition | Pros | Cons |
|---|---|---|---|
| **Multi-Prompt** | Individual prompt for each attribute and call LLM API separately. | ● More precise, relevant and comprehensive answers<br>● Error in one response doesn't affect the extraction of other attributes.<br>● Allows for tailored prompts for each attribute, thus handling complex attributes more effectively. | ● More API calls, more tokens, higher cost<br>● Sequential processing can be slower |
| **Single-Prompt** | One prompt for all attributes within a document and call API once. | ● Fewer API calls, fewer tokens, lower cost<br>● Faster especially for less complex document | ● Less precise responses<br>● Error in the response can affect the extraction of all attributes.<br>● Limited customization |

# Enhancing LLM with Retrieval Augmented Generation

**Steps for Retrieval Augmented Generation (RAG):**

1. **Prompt + Query**: The user starts the process by entering a query that reflects their question or topic of interest
2. **Query**: For each query, the system conducts a similarity search against the knowledge sources to find relevant information
3. **Enhanced Context**: Contextual information closely matching the query is extracted from the knowledge sources to serve as reference information
4. **Prompt + Query + Enhanced Text**: This extracted context, along with the user's original prompt and query, is synthesized to be inputed for the Large Language Model (LLM)
5. **Generated Text Response**: The LLM processes the integrated input to produce a precise and contextually relevant answer to the user's query

# Why we chose RAG based approach

**With RAG:**

- Can be paired with specific domain databases to provide responses in specialized fields.
- Can split longer texts into chunks.
- Can provide the source text.
- Can introduce latency.
- There's no guarantee that all information will be preserved in the result (potential information loss).

**Without RAG:**

- Works well with shorter PDFs (e.g., 1 page).
- Does not introduce latency.
- Can't work with longer PDFs (token limit issue).
- Can't locate the source text.
- Tends to hallucinate more

**We finally choose to use RAG** because:

- It allows pairing with specific domain databases for specialized responses.
- It handles longer and complex texts by splitting them into chunks.
- It provides the source text for transparency, even though there might be a latency issue.
- It offers the flexibility to work with a variety of information sources.

# Method Comparison - LLM

| Models | Pros | Cons | Decision |
|---|---|---|---|
| **GPT3.5** | ● More accurate than open-source models | ● Less accurate and advanced than GPT-4.0 | We use as final model for this project considering the balance between cost and performance |
| **GPT4.0** | ● More accurate than GPT 3.5 on complex tasks | ● More expensive than GPT 3.5 | Recommend to use for pdf reading if budget allows |
| **GPT4.0-turbo** | ● Can directly read image through API | ● New model and high cost | Recommend to use if needs to deal with highly unstructured image information |
| **Mistral-7B** | ● Completely free and secure | ● Unable to detect some information of the attributes | Not suggested |
| **Mosaic-mpt7b (instruct)** | ● Completely free and secure and accepts very long token inputs | ● Hard to convince the model to follow the exact instruction<br>● Bad performance | Not suggested |

# Result Comparison with Ground Truth
- ## Short Documents with clear tables

**Example 1:**

The GPT- 3.5 model is able to correctly locate relevant information for attributes applicable to this specific document, **achieving 100% accuracy** with no manual proofing needed for approval.
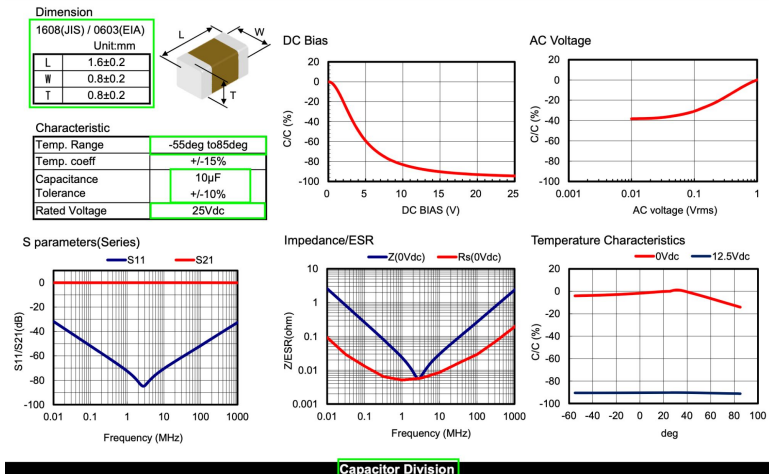


| PDF Name | Answer Type | Supplier Name | Product Type | Dimension | Voltage | Capacitance | Temperature |
|---|---|---|---|---|---|---|---|
| KGM15CR51E106K-DATA.pdf | Model Output | KYOCERA | Capacitor | 1608(JIS) / 0603(EIA) L ~ w | 25Vdc | 10uF | -55deg to85deg |
| | Ground True Value | KYOCERA | Capacitor | 1608(JIS) / 0603(EIA) L ~ w | 25V | 10 uF | -55C to 85 C |

# Result Comparison with Ground Truth

- Short Documents with clear tables

**Example 2:**

The result is still **almost 100% accurate** except for one parameter with extra width

| PDF Name | Answer Type | Supplier Name | Product Type | Dimension |
|---|---|---|---|---|
| Vzh680m2atr-1313s.pdf | Model Output | LELON ELECTRONICS | Capacitor | Length': '12.5 mm', 'Width': '10 mm', 'Height': '13.5 mm' |
| | Ground True Value | LELON ELECTRONICS | Capacitor | Length:12.5 mm, Height: 13.5mm |

| PDF Name | Answer Type | Voltage | Impedance | Capacitance | Temperature |
|---|---|---|---|---|---|
| Vzh680m2atr-1313s.pdf | Model Output | 100 V | < 0.320 (100kHz, 20°C) | 68 uF | -55°C ~ +105°C |
| | Ground True Value | 2A (100 V) | < 0.320 (100kHz, 20°C) | 68 uF | -55C to 105C |

# Result Comparison with Ground Truth
- ## Long Documents with complex tables

Example 3 document (Page 1/2)



## Example 3 (2 pages):

| PDF Name | Answer Type | Supplier Name | Product Type | Dimension |
|---|---|---|---|---|
| RK73H.pdf | Model Output | Not found | Resistor | Length: 0.016 inches (0.4 mm), Width: 0.008 inches |
| | Ground True Value | KOA Speer Electronics, Inc | Resistor | 1E (0402); 0.04 in X 0.02 in |

## Example 4 (96 pages):

| PDF Name | Answer Type | Supplier Name | Product Type | Dimension |
|---|---|---|---|---|
| IMX8MMIEC.pdf | Model Output | NXP Semiconductors | Not found | 14 x 14 mm |
| | Ground True Value | NXP Semiconductors | (App Processor) i.MX 8M M | Package Length (mm):14 mm, Package Width (mm):14 mm |

Without using the image API, PDFs which are extremely long and have complex formattings and tables, the model's accuracy is negatively impacted.

In these cases, manual check is required to validate results.

# Challenge Scenario – Future Results

- **Certain product documents requires "decoding" of part name nomenclature (see picture below)**
  - While testing, the ChatGPT API could only parse text, and therefore could not capture the formatting
  - However, the ChatGPT 4.0 web app is multimodal, and therefore can accept document images
  - NOTE: As of mid November, the OpenAI API can accommodate pictures

- **Multimodal GPT 4.0 still requires specific prompting to parse the data**
  - When asked to decode the picture (to the right), the model was generally able to summarize the tables; however there were small inaccuracies, especially with Package Type
  - When prompted to decode a part number never before seen, the model was able to properly decode the part number and list the attributes

# Challenge Scenario – Future Results



M IMX8MM @ +VT$$%A

- **Qualification level**
- **Part number series**
- **Part differentiator**

- **Silicon revision**
- **Fusing options**
- **Primary core frequency**
- **Package type – all ROHS**
- **Qualification tier Tj**

| Qualification Level | |
|---|---|
| Samples | P |
| Mass Production | M |

| Part number series | Name |
|---|---|
| IMX8MM | i.MX 8M Mini |

| Part differentiator | @ |
|---|---|
| i.MX 8M Mini Quad 4x A53, M4, GPU, VPU | 6 |
| i.MX 8M Mini QuadLite 4x A53, M4, GPU | 5 |
| i.MX 8M Mini Dual 2x A53, M4, GPU, VPU | 4 |
| i.MX 8M Mini DualLite 2x A53, M4, GPU | 3 |
| i.MX 8M Mini Solo 1x A53, M4, GPU, VPU | 2 |
| i.MX 8M Mini SoloLite 1x A53, M4, GPU | 1 |

| Temperature Tj | + |
|---|---|
| Consumer: 0 to +95℃ | D |
| Industrial: -40 to 105℃ | C |

| Package Type | ROHS |
|---|---|
| FCBGA486 14 x 14 mm, 0.5 mm pitch | VT |

| Frequency | $$ |
|---|---|
| 1.8 GHz | LZ |
| 1.6 GHz | KZ |

| Fusing | % |
|---|---|
| Default | A |
| Immersiv3D enabled w/Dolby Atmos | C |
| Immersiv3D enabled w/Dolby Atmos and DTS | D |

| Silicon rev | A |
|---|---|
| Rev A0 | A |

can you decode this

**ChatGPT**
The image you've uploaded appears to be a diagram explaining the part number nomenclature for a family of processors, specifically the i.MX 8M Mini family. The part number nomenclature is broken down into several components to specify different product options. Here's what each section seems to represent:

1. **Qualification Level**:
   - P: Samples
   - M: Mass Production
2. **Part Number Series**:
   - The name of the part number series is indicated as i.MX8M for the i.MX 8M Mini family of processors.
3. **Part Differentiator**:
   - This is a numerical value that differentiates between the types of processors within the i.MX 8M Mini family, indicating different configurations such as Quad, DualLite, Solo, etc., along with the supported features like GPU, VPU, etc.
4. **Temperature TJ**:
   - Consumer-grade temperature range: 0 to +95℃
   - Industrial-grade temperature range: -40 to +105℃
5. **Package Type**:
   - FCBGA486: This indicates the package type of the processor, which is a 486-ball Flip Chip Ball Grid Array with a certain pitch size.
6. **Silicon Revision**:
   - The letter (e.g., A) indicates the revision of the silicon used in the processor.
7. **Fusing Options**:
   - Default
   - Immersiv3D enabled w/Dolby Atmos
   - Immersiv3D enabled w/Dolby Atmos and DTS
8. **Primary Core Frequency**:
   - Different letters are used to represent the core frequency of the processor (e.g., 1.8 GHz, 1.6 GHz).
9. **Fusing, Silicon Revision, Frequency, and Package Type**:
   - Symbols (@, +, $$, %, A) and letters (VT, LZ, KZ, C, D) are used as codes to represent different options and features for the processors.

The nomenclature is essentially a structured code that allows for precise identification of processor versions and their features. It's a standardized way for manufacturers and customers to communicate about exactly which version of a processor they are referring to.

# Challenge Scenario – Future Results

**You**
using the encoding above, what is MIMX8MM5CVTKZAA

**ChatGPT**
The part number "MIMX8MM5CVTKZAA" can be decoded using the key provided in the image you've uploaded. Let's break it down:

- **MI**: This prefix is not explicitly explained in the diagram, but it typically stands for the manufacturer's code or product family.
- **MX8MM**: This refers to the product family, the i.MX 8M Mini.
- **5**: This digit is the part differentiator which specifies the type of processor. According to the key:
  - 5 indicates an i.MX 8M Mini QuadLite processor, which has 4x A53, M4, GPU cores.
- **C**: This letter indicates the temperature range. C corresponds to Industrial-grade, which is -40 to +105°C.
- **VT**: This code specifies the package type, which in this case is FCBGA486, a 486-ball Flip Chip Ball Grid Array with a 0.5 mm pitch.
- **KZ**: These two letters indicate the primary core frequency. KZ corresponds to a frequency of 1.6 GHz.
- **AA**: Typically, these letters at the end would denote the silicon revision and fusing options or additional features. However, the provided key does not detail what 'AA' stands for. In common practice, this could stand for a default configuration or a first revision without special options.

So, the part number "MIMX8MM5CVTKZAA" refers to an i.MX 8M Mini QuadLite processor with an industrial temperature range, packaged in a FCBGA486 with a 1.6 GHz core frequency, and possibly a default silicon revision and fusing options.

**You**
Do it again, but clearly denote: Qualification level, Part number series, Part differentiator, Silicon revision, Fusing options, Primary core frequency, Package type – all ROHS, Qualification tier Tj

**ChatGPT**
Certainly, let's decode the part number "MIMX8MM5CVTKZAA" with the specific sections clearly denoted:

1. **Qualification Level**:
   - **M**: Indicates that the part is approved for Mass Production.
2. **Part Number Series**:
   - **IMX8MM**: Refers to the i.MX 8M Mini series.
3. **Part Differentiator**:
   - **5**: Corresponds to the i.MX 8M Mini QuadLite version, with 4x ARM Cortex-A53 cores, 1x Cortex-M4 core, and GPU.
4. **Silicon Revision**:
   - The provided key does not specify how to decode the Silicon Revision from this part number. Normally it is indicated by a letter, but in this part number, the letter for Silicon Revision is not clear.
5. **Fusing Options**:
   - The provided key does not specify how to decode the Fusing Options from this part number. Normally it is indicated by a character, but in this part number, the character for Fusing Options is not clear.
6. **Primary Core Frequency**:
   - **KZ**: This indicates a primary core frequency of 1.6 GHz.
7. **Package Type – all ROHS**:
   - **VT**: Signifies the package type, which is FCBGA486 (14 × 14 mm, 0.5 mm pitch), and it complies with the ROHS standards for restriction of hazardous substances.
8. **Qualification Tier Tj**:
   - **C**: Denotes an industrial qualification tier with a junction temperature (Tj) range of -40 to +105°C.

In summary, the part number "MIMX8MM5CVTKZAA" describes an i.MX 8M Mini QuadLite processor, suitable for mass production, with an industrial temperature range, a 1.6 GHz primary core frequency, and packaged in a ROHS-compliant FCBGA486. The silicon revision and fusing options are not specified with the information provided.

# Demo



**Major Pipeline Steps**
- **Step 1:** PDF files are inputted and temporarily stored for preprocessing.
- **Step 2:** Textual information is extracted from the PDF files using OCR (tesseract).
- **Step 3:** The RAG (retrieval-augmented generation) module processes the text, splitting the documents into manageable segments for further analysis.
- **Step 4:** These segments are then uploaded to a vector database, where they're indexed for efficient retrieval.
- **Step 5:** A conversational chain utilizes the processed data for information extraction and question-answering (Q&A) tasks.

# Demo

# Sample LLM Costs

- **LLM Costs are calculated based on "tokens" which are groups of characters, which sometimes align with words, but not always.**
  - **1 token ~= 4 chars in English. 1 token ~= ¾ words. 100 tokens ~= 75 words**
  - Costs can be broken into Input Costs and Output Costs, both of which can be managed
  - Costs vary based on model (more powerful models are more expensive)

## Cost for 1K Documents

| Model | Context length | Cost for Input | Cost for Output | Total Costs |
|---|---|---|---|---|
| GPT 3.5-turbo | 16k | $ 3.80 | $ 1.00 | $ 4.80 |
| GPT 4.0-turbo | 128k | $ 38.00 | $ 15.00 | $ 53.00 |
| GPT 4.0 | 32k | $ 228.00 | $ 60.00 | $ 288.00 |

Note: costs assumes average input of 3800 tokens and output of 200 tokens
OpenAI Pricing: https://openai.com/pricing

# Next Steps and Recommendations

1. **OpenAI just released picture handling in their new API; based on initial testing, we think this method will be immensely useful for Cisco:**
   a. Product documents that require "decoding" as previously mentioned
   b. Reading and summarizing graphs and line charts
   c. Reading and interpreting data tables (initial tests conversion of data tables to JSON format were successful, though will require prompt tuning)
2. **Human-in-the-loop validation: complex document types and documents of a certain length would use humans for validation while simple, shorter documents would rely solely on the pipeline**
3. **Customized prompt with attribute explanation works better than a generic prompt, but prompt engineering effort should be considered.**
4. **Given the additional model performance and extended context window, we suggest future builds use the GPT4.0-turbo model**

# Next Steps and Recommendations (cont.)

- To determine the final prompt of choice, we suggest following our defined "Success Framework" which scores Model Output against Ground Truth across many documents (>30 preferred)
- Given the constantly changing technical environment, we suggest building the final pipeline in a modular fashion by keeping the RAG, Model Itself, File Reading Mechanism, and Web App independent of one another to allow for simple "lift and shifts"

# Milestones (Proposed)

1. **Stage 1: Single Document Attribute Extraction (Oct 1st - Oct 15)**
   - Use a variety of methods; results will dictate which method(s) we'll continue forward with
2. **Stage 2: Cross-Product Attribute Extraction from Same Supplier (Oct 16 - Oct 31)**
   - Continue to narrow down on methodology
3. **Mid-Stage Review (Nov 1st)**

**We are here**

---

4. **Stage 3: Cross-Supplier and Cross-Product Attribute Extraction (Nov 1st - Nov 15)**
   - Test the extraction process on suppliers and product data, refine the algorithm based on test results
   - **(updated)** Validate the accuracy, consistency, and costs of the final attribute extraction method
   - **(new)** Get approval from the Cisco team on the final architecture
5. **(updated) Stage 4: Chat-Bot, Web Application, and Final Pipeline Development (Nov 16 - Dec 6th)**
6. **Final Review and Presentation (Dec 6th)**
   - Review the entire project and the milestones achieved, gather feedback for future improvements and iterations

# Thank You!

# Questions?

# Appendix

- **Project Proposal Link**
- **Project Proposal Slides**
- **Process Tracker**

# Sample Tesseract Output

- **Step 1:  Conduct a Manual Inspection of Attributes in the Selected PDFs**
  - The MLDS team will meticulously review the final 30 PDF documents to identify the most accurate ground truth values for each attribute. This will serve as our reference benchmark for subsequent validation checks

- **Step 2:  Compare the Model Output with Ground Truth**
  - In-depth grading rubric will be developed to determine partial accuracy scores

- **Step 3: Quantify Model Accuracy Across Documents**
  - To determine the model's accuracy, divide the number of attributes the model correctly identified by the total number of attributes across all documents

# Success Framework

- **Step 1:  Conduct a Manual Inspection of Attributes in the Selected PDFs**
    - The MLDS team will meticulously review the final 30 PDF documents to identify the most accurate ground truth values for each attribute. This will serve as our reference benchmark for subsequent validation checks

- **Step 2:  Compare the Model Output with Ground Truth**
    - In-depth grading rubric will be developed to determine partial accuracy scores

- **Step 3: Quantify Model Accuracy Across Documents**
    - To determine the model's accuracy, divide the number of attributes the model correctly identified by the total number of attributes across all documents

# Final Deliverables

- **Deck on Document Extraction: Technology Comparison & Framework**
  - Pros and cons comparison of existing technologies
  - Suggested document extraction framework
  - Potential next steps
- **LLM Prompt Engineering**
  - Explore best system message, questions format, and LLM parameters (in progress)
- **WebApp Link & Demo**
  - Interactive conversational UI with the ability to chat with an uploaded documents
  - Run document intelligence AI across a collection of documents and store extracted features in structured format.

# Potential Blockers

- **Tesseract and GPT 4.0 approvals**
    - Internal Cisco review process may dictate which solutions are viable
- **LLM Costs**
    - Costs could balloon with the addition of document chat functionality, by-attribute querying, and size of the final dataset (ex. applying the pipeline to thousands of documents)

### GPT API Costs

| Model | Input | Output |
| --- | --- | --- |
| GPT 3.5 Turbo - 4k context | $0.0015 per 1K tokens | $0.002 per 1K tokens |
| GPT 3.5 Turbo - 16k context | $0.003 per 1K tokens | $0.004 per 1K tokens |
| GPT 4.0 - 8K context | $0.03 per 1K tokens | $0.06 per 1K tokens |
| GPT 4.0 - 32K context | $0.06 per 1K tokens | $0.12 per 1K tokens |

Northwestern

# Proposed Methods (From Midpoint)

Based on our testing , here is our proposed methodology:

- **Attribute Extraction Pipeline:** extracts a set list of attributes along with their context/proof in JSON format. Attributes are then added to an existing database or filestore system
  - Solutions: GPT 3.5, GPT 4.0 with or without RAG, running on an EC2 or ECS instance
- **Storage**: stores product attributes and their context for future querying
  - Solutions: Relational database or object storage (ex. S3 bucket)
- **Web App**: allows end users to upload documents or directory and extract information from the files.
  - Solutions: Streamlit
- **Chatbot**: allows users to ask questions about content from existing documents or submit new documents for querying purposes.
  - Solutions: LLM + Streamlit

# Method Comparison - Example

**One prompt for all attributes:**

Answer: Supplier Name: AICC (Assumed based on the document header)

Product Name: EXAMAX2™ and EXAMEZZ2™ Connector System

Product Type: Connector System

Dimensions: The document does not explicitly state the dimensions of the connector system.

Orientation: The document mentions various orientations for different configurations, such as vertical header (VH), right angle receptacle (RAR), right angle orthogonal header (RAOH), and right angle header (RAH). Specific orientations may vary depending on the configuration.

Current Rating: The document refers to AICC Power Design Guide GS-20-0456 for the current rating. This guide should provide the relevant information.

Voltage: The document specifies the operating voltage rating as follows:
- Agency Voltage Rating < 30 VAC RMS / DC
- Non-Agency Voltage Rating = 150 VAC RMS or VDC Maximum

Frequency: The document does not mention a specific frequency.

Impedance / Capacitance: The document does not specify impedance or capacitance values.

Temperatures: The document provides the following temperature ranges:
- Operating Temperature Range = -55 to 85 degrees Celsius
- High Temperature Life Test Temperature = 85 degrees Celsius
- Thermal Shock Test Temperature Range = -55 to +85 degrees Celsius
- Humidity Test Temperature Range = 25 to 65 degrees Celsius

Please note that for specific technical details and performance requirements, it is recommended to consult the AICC Power Design Guide GS-20-0456 or other relevant documentation provided by the supplier.

# Method Comparison - Mosaic Example

| Document Name | KGM15CR51E106K-DATA.pdf | gs-12-1565.pdf |
|---|---|---|
| Document Length | 1 page long | 12 page long |
| Model Output | The electrical characteristics data for this capacitor are as follows:<br>**Supplier name:** KEMET<br>**Part number:** CM105X5R106K25A<br>**Product type:** Ceramic Capacitors<br>**Case size:** 0603<br>**Lead spacing:** X5R<br>**Dimensions:** 1608(JIS)/0603(EIA) | Chained Grit Channels Chose Extensive EXCLUSIVE EX Chasers Chancing Chract Chrut Chrst Exact Chasing Exposed Ch Chrit Ch Ch Ch Ch Expose Ch Ch Exciting Ch Ex Ch Ch Ch Ch Ch Ch Exposing Ch Ch Ex Ch Ch Ch Ex Ch Ch Ch Ch |
| Response Time | 3 seconds | 3 minutes |

# Prompt Engineering [other prompts we tested]

| Sample Prompts | Observations |
|---|---|
| "I'm looking for information in this document. Please extract the following information: You can **either provide the extracted information or provide the range of values** for the following: (a) Supplier Name, (b) Product type, (c) Dimensions, (d) Orientation if any Current Rating, (e) Voltage, (f) Frequency, (g) Impedance, (h) Capacitance, (i) Temperatures, (j) Other configuration settings (all the information that is not covered in the above categories)" | ● Output **texts** with a list of attributes and results. <br> ● Able to find and provide some additional attributes information with a range of value (if provided in pdf) compared with the previous one. |
| System message: Use the following pieces of context to answer the question at the end. If you don't know the answer, just say that you don't know, don't try to make up an answer. <br> Prompt: Please **extract any sentence** about the following information and include note for each attribute you extracted. (a) supplier name,..., (j) Other configuration settings. Please provide the answer **in the following format:** { "attribute": "attribute Name", "value": "attribute value", "note": "context to where and why you extracted the value" }, { "attribute": "attribute Name", "value": "attribute value", "context": "context to where and why you extracted the value" }, etc. for Supplier Name, Product type, Dimensions and so on. | ● Output **JSON** with dictionaries of attribute names, attribute values and contexts. <br> ● GPT 4 is able to output in the instructed format, and provide context for where its answers come from (either section number or section title) <br> ● e.g.{ "attribute": "Frequency", "value": "10 to 500 to 10 hertz", "note": "This value was extracted from the section 8.6 Vibration Sinusoidal – EIA 364 -28 where it specifies the Frequency Range for the vibration test." } |

# Method Comparison - LLM

| Models | Pros | Cons | Decision |
|---|---|---|---|
| **GPT3.5** | ● More accurate than open-source models | ● Less accurate and advanced than GPT-4.0 | We use as final model for this project considering the balance between cost and performance |
| **GPT4.0** | ● More accurate than GPT 3.5 on complex tasks | ● More expensive than GPT 3.5 | Recommend to use for pdf reading if budget allows |
| **GPT4.0-turbo** | ● Can directly read image through API | ● New model and high cost | Recommend to use if needs to deal with highly unstructured image information |
| **Mistral-7B** | ● Completely free and secure | ● Unable to detect some information of the attributes | Not suggested |
| **Mosaic-mpt7b (instruct)** | ● Completely free and secure and accepts very long token inputs | ● Hard to convince the model to follow the exact instruction<br>● Bad performance | Not suggested |