

# Pointing to External Data and Custom Parsing

In DocETL, you have full control over your dataset JSONs. These JSONs typically contain objects with key-value pairs, where you can reference external files that you want to process in your pipeline. This referencing mechanism, which we call "pointing", allows DocETL to locate and process external files that require special handling before they can be used in your main pipeline.



## Why Use Custom Parsing?

Consider these scenarios where custom parsing of referenced files is beneficial:

- Your dataset JSON references Excel spreadsheets containing sales data.
- You have entries pointing to scanned receipts in PDF format that need OCR processing.
- You want to extract text from Word documents or PowerPoint presentations by referencing their file locations.

In these cases, custom parsing enables you to transform your raw external data into a format that DocETL can process effectively within your pipeline. The pointing mechanism allows DocETL to locate these external files and apply custom parsing seamlessly. *(Pointing in DocETL refers to the practice of including references or paths to external files within your dataset JSON. Instead of embedding the entire content of these files, you simply "point" to their locations, allowing DocETL to access and process them as needed during the pipeline execution.)*

## Dataset JSON Example

Let's look at a typical dataset JSON file that you might create:

```
[
  { "id": 1, "excel_path": "sales_data/january_sales.xlsx" },
  { "id": 2, "excel_path": "sales_data/february_sales.xlsx" }
]
```

In this example, you've specified paths to Excel files. DocETL will use these paths to locate and process the external files. However, without custom parsing, DocETL wouldn't know how to handle the contents of these files. This is where parsing tools come in handy.

# Custom Parsing in Action

## 1. Configuration

To use custom parsing, you need to define parsing tools in your DocETL configuration file. Here's an example:

```
parsing_tools:
- name: top_products_report
  function_code: |
    def top_products_report(document: Dict) -> List[Dict]:
        import pandas as pd

        # Read the Excel file
        filename = document["excel_path"]
        df = pd.read_excel(filename)

        # Calculate total sales
        total_sales = df['Sales'].sum()

        # Find top 500 products by sales
        top_products = df.groupby('Product')['Sales'].sum().nlargest(500)

        # Calculate month-over-month growth
        df['Date'] = pd.to_datetime(df['Date'])
        monthly_sales = df.groupby(df['Date'].dt.to_period('M'))
['Sales'].sum()
        mom_growth = monthly_sales.pct_change().fillna(0)

        # Prepare the analysis report
        report = [
            f"Total Sales: ${total_sales:,.2f}",
            "\nTop 500 Products by Sales:",
            top_products.to_string(),
            "\nMonth-over-Month Growth:",
            mom_growth.to_string()
        ]

        # Return a list of dicts representing the output
        # The input document will be merged into each output doc,
        # so we can access all original fields from the input doc.
        return [{"sales_analysis": "\n".join(report)}]

datasets:
  sales_reports:
    type: file
    source: local
    path: "sales_data/sales_paths.json"
    parsing:
      - function: top_products_report

  receipts:
    type: file
    source: local
    path: "receipts/receipt_paths.json"
```

```
parsing:
  - input_key: pdf_path
    function: paddleocr_pdf_to_string
    output_key: receipt_text
    ocr_enabled: true
    lang: "en"
```

In this configuration:

- We define a custom `top_products_report` function for Excel files.
- We use the built-in `paddleocr_pdf_to_string` parser for PDF files.
- We apply these parsing tools to the external files referenced in the respective datasets.

## 2. Pipeline Integration

Once you've defined your parsing tools and datasets, you can use the processed data in your pipeline:

```
pipeline:
  steps:
    - name: process_sales
      input: sales_reports
      operations:
        - summarize_sales
    - name: process_receipts
      input: receipts
      operations:
        - extract_receipt_info
```

This pipeline will use the parsed data from both Excel files and PDFs for further processing.

## How Data Gets Parsed and Formatted

When you run your DocETL pipeline, the parsing tools you've specified in your configuration file are applied to the external files referenced in your dataset JSONs.

Here's what happens:

1. DocETL reads your dataset JSON file.
2. For each entry in the dataset, it looks at the parsing configuration you've specified.
3. It applies the appropriate parsing function to the file path provided in the dataset JSON.
4. The parsing function processes the file and returns the data in a format DocETL can work with (typically a list of strings).

Let's look at how this works for our earlier examples:

## Excel Files (using top\_products\_report)

For an Excel file like "sales\_data/january\_sales.xlsx":

- The top\_products\_report function reads the Excel file.
- It processes the sales data and generates a report of top-selling products.
- The output might look like this:

```
Top Products Report - January 2023

1. Widget A - 1500 units sold
2. Gadget B - 1200 units sold
3. Gizmo C - 950 units sold
4. Doohickey D - 800 units sold
5. Thingamajig E - 650 units sold
...

Total Revenue: $245,000
Best Selling Category: Electronics
```

## PDF Files (using paddleocr\_pdf\_to\_string)

For a PDF file like "receipts/receipt001.pdf":

- The paddleocr\_pdf\_to\_string function reads the PDF file.
- It uses PaddleOCR to perform optical character recognition on each page.
- The function combines the extracted text from all pages into a single string. The output might look like this:

```
RECEIPT
Store: Example Store
Date: 2023-05-15
Items:

1. Product A - $10.99
2. Product B - $15.50
3. Product C - $7.25
4. Product D - $22.00
  Subtotal: $55.74
  Tax (8%): $4.46
  Total: $60.20

Payment Method: Credit Card
Card Number: \*\*\* \* \*\*\* \*\*\*\* 1234

Thank you for your purchase!
```

This parsed and formatted data is then passed to the respective operations in your pipeline for further processing.

## Running the Pipeline

Once you've set up your pipeline configuration file with the appropriate parsing tools and dataset definitions, you can run your DocETL pipeline. Here's how:

1. Ensure you have DocETL installed in your environment.
2. Open a terminal or command prompt.
3. Navigate to the directory containing your pipeline configuration file.
4. Run the following command:

```
docetl run pipeline.yaml
```

Replace `pipeline.yaml` with the name of your pipeline file if it's different.

When you run this command:

1. DocETL reads your pipeline file.
2. It processes each dataset using the specified parsing tools.
3. The pipeline steps are executed in the order you defined.
4. Any operations you've specified (like `summarize_sales` or `extract_receipt_info`) are applied to the parsed data.
5. The results are saved according to your output configuration.

## Built-in Parsing Tools

DocETL provides several built-in parsing tools to handle common file formats and data processing tasks. These tools can be used directly in your configuration by specifying their names in the `function` field of your parsing tools configuration. Here's an overview of the available built-in parsing tools:

Convert an Excel file to a string representation or a list of string representations.

### Parameters:

Name	Type	Description	Default
<code>filename</code>	<code>str</code>	Path to the <code>xlsx</code> file.	<i>required</i>
<code>orientation</code>	<code>str</code>	Either "row" or "col" for cell arrangement.	<code>'col'</code>

Name	Type	Description	Default
<code>col_order</code>	<code>list[str]   None</code>	List of column names to specify the order.	<code>None</code>
<code>doc_per_sheet</code>	<code>bool</code>	If True, return a list of strings, one per sheet.	<code>False</code>

**Returns:**

Type	Description
<code>list[str]</code>	<code>list[str]</code> : String representation(s) of the Excel file content.

## Source code in docetl/parsing\_tools.py

```

99 @with_input_output_key
100 def xlsx_to_string(
101     filename: str,
102     orientation: str = "col",
103     col_order: list[str] | None = None,
104     doc_per_sheet: bool = False,
105 ) -> list[str]:
106     """
107     Convert an Excel file to a string representation or a list of string
108     representations.
109
110     Args:
111         filename (str): Path to the xlsx file.
112         orientation (str): Either "row" or "col" for cell arrangement.
113         col_order (list[str] | None): List of column names to specify the
114         order.
115         doc_per_sheet (bool): If True, return a list of strings, one per
116         sheet.
117
118     Returns:
119         list[str]: String representation(s) of the Excel file content.
120     """
121     import openpyxl
122
123     wb = openpyxl.load_workbook(filename)
124
125     def process_sheet(sheet):
126         if col_order:
127             headers = [
128                 col for col in col_order if col in sheet.iter_cols(1,
129 sheet.max_column)
130             ]
131         else:
132             headers = [cell.value for cell in sheet[1]]
133
134         result = []
135         if orientation == "col":
136             for col_idx, header in enumerate(headers, start=1):
137                 column = sheet.cell(1, col_idx).column_letter
138                 column_values = [cell.value for cell in sheet[column]
139 [1:]]
140                 result.append(f"{header}: " + "\n".join(map(str,
141 column_values)))
142                 result.append("") # Empty line between columns
143             else: # row
144                 for row in sheet.iter_rows(min_row=2, values_only=True):
145                     row_dict = {
146                         header: value for header, value in zip(headers, row)
147                     }
148                     if header:
149                         result.append(
150                             " | ".join(
151                                 [f"{header}: {value}" for header, value in
152 row_dict.items()]
153                             )
154                         )

```

```
        return "\n".join(result)

    if doc_per_sheet:
        return [process_sheet(sheet) for sheet in wb.worksheets]
    else:
        return [process_sheet(wb.active)]
```

options: show\_root\_heading: true heading\_level: 3

Read the content of a text file and return it as a list of strings (only one element).

Parameters:

Name	Type	Description	Default
filename	str	Path to the txt or md file.	required

Returns:

Type	Description
list[str]	list[str]: Content of the file as a list of strings.

Source code in docetl/parsing\_tools.py

```
156 @with_input_output_key
157 def txt_to_string(filename: str) -> list[str]:
158     """
159     Read the content of a text file and return it as a list of strings
160     (only one element).
161
162     Args:
163         filename (str): Path to the txt or md file.
164
165     Returns:
166         list[str]: Content of the file as a list of strings.
167     """
168     with open(filename, "r", encoding="utf-8") as file:
169         return [file.read()]
```

options: show\_root\_heading: true heading\_level: 3

Extract text from a Word document.

Parameters:



Name	Type	Description	Default
<code>filename</code>	<code>str</code>	Path to the docx file.	<i>required</i>

**Returns:**

Type	Description
<code>list[str]</code>	<code>list[str]</code> : Extracted text from the document.

Source code in `docetl/parsing_tools.py`

```
171 @with_input_output_key
172 def docx_to_string(filename: str) -> list[str]:
173     """
174     Extract text from a Word document.
175
176     Args:
177         filename (str): Path to the docx file.
178
179     Returns:
180         list[str]: Extracted text from the document.
181     """
182     from docx import Document
183
184     doc = Document(filename)
185     return ["\n".join([paragraph.text for paragraph in doc.paragraphs])]
```

options: show\_root\_heading: true heading\_level: 3

Transcribe speech from an audio file to text using Whisper model via litellm. If the file is larger than 25 MB, it's split into 10-minute chunks with 30-second overlap.

**Parameters:**

Name	Type	Description	Default
<code>filename</code>	<code>str</code>	Path to the mp3 or mp4 file.	<i>required</i>

**Returns:**

Type	Description
<code>list[str]</code>	<code>list[str]</code> : Transcribed text.

#### Source code in `docetl/parsing_tools.py`

```

52 @with_input_output_key
53 def whisper_speech_to_text(filename: str) -> list[str]:
54     """
55     Transcribe speech from an audio file to text using Whisper model via
56     litellm.
57     If the file is larger than 25 MB, it's split into 10-minute chunks
58     with 30-second overlap.
59
60     Args:
61         filename (str): Path to the mp3 or mp4 file.
62
63     Returns:
64         list[str]: Transcribed text.
65     """
66
67     from litellm import transcription
68
69     file_size = os.path.getsize(filename)
70     if file_size > 25 * 1024 * 1024: # 25 MB in bytes
71         from pydub import AudioSegment
72
73         audio = AudioSegment.from_file(filename)
74         chunk_length = 10 * 60 * 1000 # 10 minutes in milliseconds
75         overlap = 30 * 1000 # 30 seconds in milliseconds
76
77         chunks = []
78         for i in range(0, len(audio), chunk_length - overlap):
79             chunk = audio[i : i + chunk_length]
80             chunks.append(chunk)
81
82         transcriptions = []
83
84         for i, chunk in enumerate(chunks):
85             buffer = io.BytesIO()
86             buffer.name = f"temp_chunk_{i}_{os.path.basename(filename)}"
87             chunk.export(buffer, format="mp3")
88             buffer.seek(0) # Reset buffer position to the beginning
89
90             response = transcription(model="whisper-1", file=buffer)
91             transcriptions.append(response.text)
92
93         return transcriptions
94     else:
95         with open(filename, "rb") as audio_file:
96             response = transcription(model="whisper-1", file=audio_file)
97
98         return [response.text]

```

options: show\_root\_heading: true heading\_level: 3

Extract text from a PowerPoint presentation.

#### Parameters:

Name	Type	Description	Default
<code>filename</code>	<code>str</code>	Path to the pptx file.	<i>required</i>
<code>doc_per_slide</code>	<code>bool</code>	If True, return each slide as a separate document. If False, return the entire presentation as one document.	<code>False</code>

#### Returns:

Type	Description
<code>list[str]</code>	<code>list[str]</code> : Extracted text from the presentation. If <code>doc_per_slide</code> is True, each string in the list represents a single slide. Otherwise, the list contains a single string with all slides' content.

Source code in `docetl/parsing_tools.py`

```

188 @with_input_output_key
189 def pptx_to_string(filename: str, doc_per_slide: bool = False) ->
190     list[str]:
191         """
192         Extract text from a PowerPoint presentation.
193
194         Args:
195             filename (str): Path to the pptx file.
196             doc_per_slide (bool): If True, return each slide as a separate
197             document. If False, return the entire presentation as one
198             document.
199
200         Returns:
201             list[str]: Extracted text from the presentation. If doc_per_slide
202             is True, each string in the list represents a single slide.
203             Otherwise, the list contains a single string with all slides'
204             content.
205         """
206         from pptx import Presentation
207
208         prs = Presentation(filename)
209         result = []
210
211         for slide in prs.slides:
212             slide_content = []
213             for shape in slide.shapes:
214                 if hasattr(shape, "text"):
215                     slide_content.append(shape.text)
216
217             if doc_per_slide:
218                 result.append("\n".join(slide_content))
219             else:
220                 result.extend(slide_content)
221
222         if not doc_per_slide:
223             result = ["\n".join(result)]
224
225         return result

```

options: show\_root\_heading: true heading\_level: 3

Note to developers: We used [this documentation](#) as a reference.

This function uses Azure Document Intelligence to extract text from documents. To use this function, you need to set up an Azure Document Intelligence resource:

1. [Create an Azure account](#) if you don't have one
2. Set up a Document Intelligence resource in the [Azure portal](#)
3. Once created, find the resource's endpoint and key in the Azure portal
4. Set these as environment variables:

5. DOCUMENTINTELLIGENCE\_API\_KEY: Your Azure Document Intelligence API key

6. DOCUMENTINTELLIGENCE\_ENDPOINT: Your Azure Document Intelligence endpoint URL

The function will use these credentials to authenticate with the Azure service. If the environment variables are not set, the function will raise a ValueError.

The Azure Document Intelligence client is then initialized with these credentials. It sends the document (either as a file or URL) to Azure for analysis. The service processes the document and returns structured information about its content.

This function then extracts the text content from the returned data, applying any specified formatting options (like including line numbers or font styles). The extracted text is returned as a list of strings, with each string representing either a page (if doc\_per\_page is True) or the entire document.

#### Parameters:

Name	Type	Description	Default
<code>filename</code>	<code>str</code>	Path to the file to be analyzed or URL of the document if use_url is True.	<i>required</i>
<code>use_url</code>	<code>bool</code>	If True, treat filename as a URL. Defaults to False.	<code>False</code>
<code>include_line_numbers</code>	<code>bool</code>	If True, include line numbers in the output. Defaults to False.	<code>False</code>
<code>include_handwritten</code>	<code>bool</code>	If True, include handwritten text in the output. Defaults to False.	<code>False</code>
<code>include_font_styles</code>	<code>bool</code>	If True, include font style information in the output. Defaults to False.	<code>False</code>
<code>include_selection_marks</code>	<code>bool</code>	If True, include selection marks in the output.	<code>False</code>

Name	Type	Description	Default
		Defaults to False.	
<code>doc_per_page</code>	<code>bool</code>	If True, return each page as a separate document. Defaults to False.	<code>False</code>

**Returns:**

Type	Description
<code>list[str]</code>	<code>list[str]</code> : Extracted text from the document. If <code>doc_per_page</code> is True, each string in the list represents a single page. Otherwise, the list contains a single string with all pages' content.

**Raises:**

Type	Description
<code>ValueError</code>	If <code>DOCUMENTINTELLIGENCE_API_KEY</code> or <code>DOCUMENTINTELLIGENCE_ENDPOINT</code> environment variables are not set.

**Source code in `docetl/parsing_tools.py`**

```

226 @with_input_output_key
227 def azure_di_read(
228     filename: str,
229     use_url: bool = False,
230     include_line_numbers: bool = False,
231     include_handwritten: bool = False,
232     include_font_styles: bool = False,
233     include_selection_marks: bool = False,
234     doc_per_page: bool = False,
235 ) -> list[str]:
236     """
237     > Note to developers: We used [this documentation]
238     (https://learn.microsoft.com/en-us/azure/ai-services/document-
239     intelligence/how-to-guides/use-sdk-rest-api?view=doc-intel-
240     4.0.0&tabs=windows&pivots=programming-language-python) as a reference.
241
242     This function uses Azure Document Intelligence to extract text from
243     documents.
244
245     To use this function, you need to set up an Azure Document
246     Intelligence resource:
247
248     1. [Create an Azure account](https://azure.microsoft.com/) if you
249     don't have one
250     2. Set up a Document Intelligence resource in the [Azure portal]
251     (https://portal.azure.com/#create/Microsoft.CognitiveServicesFormRecogniz
252     er)
253     3. Once created, find the resource's endpoint and key in the Azure
254     portal
255     4. Set these as environment variables:
256         - DOCUMENTINTELLIGENCE_API_KEY: Your Azure Document Intelligence
257         API key
258         - DOCUMENTINTELLIGENCE_ENDPOINT: Your Azure Document Intelligence
259         endpoint URL
260
261     The function will use these credentials to authenticate with the
262     Azure service.
263
264     If the environment variables are not set, the function will raise a
265     ValueError.
266
267     The Azure Document Intelligence client is then initialized with these
268     credentials.
269
270     It sends the document (either as a file or URL) to Azure for
271     analysis.
272
273     The service processes the document and returns structured information
274     about its content.
275
276     This function then extracts the text content from the returned data,
277     applying any specified formatting options (like including line
278     numbers or font styles).
279
280     The extracted text is returned as a list of strings, with each string
281     representing either a page (if doc_per_page is True) or the entire
282     document.
283
284     Args:
285         filename (str): Path to the file to be analyzed or URL of the
286         document if use_url is True.
287         use_url (bool, optional): If True, treat filename as a URL.

```

```

283 Defaults to False.
284     include_line_numbers (bool, optional): If True, include line
285 numbers in the output. Defaults to False.
286     include_handwritten (bool, optional): If True, include
287 handwritten text in the output. Defaults to False.
288     include_font_styles (bool, optional): If True, include font style
289 information in the output. Defaults to False.
290     include_selection_marks (bool, optional): If True, include
291 selection marks in the output. Defaults to False.
292     doc_per_page (bool, optional): If True, return each page as a
293 separate document. Defaults to False.
294
295 Returns:
296     list[str]: Extracted text from the document. If doc_per_page is
297 True, each string in the list represents
298         a single page. Otherwise, the list contains a single
299 string with all pages' content.
300
301 Raises:
302     ValueError: If DOCUMENTINTELLIGENCE_API_KEY or
303 DOCUMENTINTELLIGENCE_ENDPOINT environment variables are not set.
304 """
305
306 from azure.ai.documentintelligence import DocumentIntelligenceClient
307 from azure.ai.documentintelligence.models import
308 AnalyzeDocumentRequest
309 from azure.core.credentials import AzureKeyCredential
310
311 key = os.getenv("DOCUMENTINTELLIGENCE_API_KEY")
312 endpoint = os.getenv("DOCUMENTINTELLIGENCE_ENDPOINT")
313
314 if key is None:
315     raise ValueError("DOCUMENTINTELLIGENCE_API_KEY environment
316 variable is not set")
317 if endpoint is None:
318     raise ValueError(
319         "DOCUMENTINTELLIGENCE_ENDPOINT environment variable is not
320 set"
321     )
322
323 document_analysis_client = DocumentIntelligenceClient(
324     endpoint=endpoint, credential=AzureKeyCredential(key)
325 )
326
327 if use_url:
328     poller = document_analysis_client.begin_analyze_document(
329         "prebuilt-read", AnalyzeDocumentRequest(url_source=filename)
330     )
331 else:
332     with open(filename, "rb") as f:
333         poller =
334 document_analysis_client.begin_analyze_document("prebuilt-read", f)
335
336 result = poller.result()
337
338 style_content = []
339 content = []
340
341 if result.styles:
342     for style in result.styles:
343         if style.is_handwritten and include_handwritten:

```



```

344         handwritten_text = ",".join(
345             [
346                 result.content[span.offset : span.offset +
347 span.length]
348                 for span in style.spans
349             ]
350         )
351         style_content.append(f"Handwritten content:
352 {handwritten_text}")
353
354         if style.font_style and include_font_styles:
355             styled_text = ",".join(
356                 [
357                     result.content[span.offset : span.offset +
358 span.length]
359                     for span in style.spans
360                 ]
361             )
362             style_content.append(f"'{style.font_style}' font style:
{styled_text}")

    for page in result.pages:
        page_content = []

        if page.lines:
            for line_idx, line in enumerate(page.lines):
                if include_line_numbers:
                    page_content.append(f" Line #{line_idx}:
{line.content}")
                else:
                    page_content.append(f"{line.content}")

            if page.selection_marks and include_selection_marks:
                # TODO: figure this out
                for selection_mark_idx, selection_mark in
enumerate(page.selection_marks):
                    page_content.append(
                        f"Selection mark #{selection_mark_idx}: State is
'{selection_mark.state}' within bounding polygon "
                        f"'{selection_mark.polygon}' and has a confidence of
{selection_mark.confidence}"
                    )

            content.append("\n".join(page_content))

    if doc_per_page:
        return style_content + content
    else:
        return [
            "\n\n".join(
                [
                    "\n".join(style_content),
                    "\n\n".join(
                        f"Page {i+1}: \n{page_content}"
                        for i, page_content in enumerate(content)
                    ),
                ]
            ),
        ]

```

options: heading\_level: 3 show\_root\_heading: true

Extract text and image information from a PDF file using PaddleOCR for image-based PDFs.

**Note: this is very slow!!**

**Parameters:**

Name	Type	Description	Default
<code>input_path</code>	<code>str</code>	Path to the input PDF file.	<i>required</i>
<code>doc_per_page</code>	<code>bool</code>	If True, return a list of strings, one per page. If False, return a single string.	<code>False</code>
<code>ocr_enabled</code>	<code>bool</code>	Whether to enable OCR for image-based PDFs.	<code>True</code>
<code>lang</code>	<code>str</code>	Language of the PDF file.	<code>'en'</code>

**Returns:**

Type	Description
<code>list[str]</code>	<code>list[str]</code> : Extracted content as a list of formatted strings.

Source code in `docetl/parsing_tools.py`

```

365 @with_input_output_key
366 def paddleocr_pdf_to_string(
367     input_path: str,
368     doc_per_page: bool = False,
369     ocr_enabled: bool = True,
370     lang: str = "en",
371 ) -> list[str]:
372     """
373     Extract text and image information from a PDF file using PaddleOCR
374     for image-based PDFs.
375
376     **Note: this is very slow!!**
377
378     Args:
379         input_path (str): Path to the input PDF file.
380         doc_per_page (bool): If True, return a list of strings, one per
381         page.
382             If False, return a single string.
383         ocr_enabled (bool): Whether to enable OCR for image-based PDFs.
384         lang (str): Language of the PDF file.
385
386     Returns:
387         list[str]: Extracted content as a list of formatted strings.
388     """
389     import fitz
390     import numpy as np
391     from paddleocr import PaddleOCR
392
393     ocr = PaddleOCR(use_angle_cls=True, lang=lang)
394
395     pdf_content = []
396
397     with fitz.open(input_path) as pdf:
398         for page_num in range(len(pdf)):
399             page = pdf[page_num]
400             text = page.get_text()
401             images = []
402
403             # Extract image information
404             for img_index, img in enumerate(page.get_images(full=True)):
405                 rect = page.get_image_bbox(img)
406                 images.append(f"Image {img_index + 1}: bbox {rect}")
407
408             page_content = f"Page {page_num + 1}:\n"
409             page_content += f"Text:\n{text}\n"
410             page_content += f"Images:\n" + "\n".join(images) + "\n"
411
412             if not text and ocr_enabled:
413                 mat = fitz.Matrix(2, 2)
414                 pix = page.get_pixmap(matrix=mat)
415                 img = np.frombuffer(pix.samples, dtype=np.uint8).reshape(
416                     pix.height, pix.width, 3
417                 )
418
419                 ocr_result = ocr.ocr(img, cls=True)
420                 page_content += f"OCR Results:\n"
421                 for line in ocr_result[0]:

```

```
422         bbox, (text, _) = line
423         page_content += f"{bbox}, {text}\n"
424
425     pdf_content.append(page_content)
426
427     if not doc_per_page:
428         return ["\n\n".join(pdf_content)]
429
430     return pdf_content
```

options: heading\_level: 3 show\_root\_heading: true

## Using Function Arguments with Parsing Tools

When using parsing tools in your DocETL configuration, you can pass additional arguments to the parsing functions.

For example, when using the `xlsx_to_string` parsing tool, you can specify options like the orientation of the data, the order of columns, or whether to process each sheet separately. Here's an example of how to use such kwargs in your configuration:

```
datasets:
  my_sales:
    type: file
    source: local
    path: "sales_data/sales_paths.json"
    parsing_tools:
      - name: excel_parser
        function: xlsx_to_string
        orientation: row
        col_order: ["Date", "Product", "Quantity", "Price"]
        doc_per_sheet: true
```

## Contributing Built-in Parsing Tools

While DocETL provides several built-in parsing tools, the community can always benefit from additional utilities. If you've developed a parsing tool that you think could be useful for others, consider contributing it to the DocETL repository. Here's how you can add new built-in parsing utilities:

1. Fork the DocETL repository on GitHub.
2. Clone your forked repository to your local machine.
3. Navigate to the `docetl/parsing_tools.py` file.
4. Add your new parsing function to this file. The function should also be added to the `PARSING_TOOLS` dictionary.
5. Update the documentation in the function's docstring.

6. Create a pull request to merge your changes into the main DocETL repository.



#### Guidelines for Contributing Parsing Tools

When contributing a new parsing tool, make sure it follows these guidelines:

- The function should have a clear, descriptive name.
- Include comprehensive docstrings explaining the function's purpose, parameters, and return value. The return value should be a list of strings.
- Handle potential errors gracefully and provide informative error messages.
- If your parser requires additional dependencies, make sure to mention them in the pull request.

## Creating Custom Parsing Tools

If the built-in tools don't meet your needs, you can create your own custom parsing tools. Here's how:

1. Define your parsing function in the `parsing_tools` section of your configuration.
2. Ensure your function takes a item (dict) as input and returns a list of items (dicts).
3. Use your custom parser in the `parsing` section of your dataset configuration.

For example:

```
parsing_tools:
- name: my_custom_parser
  function_code: |
    def my_custom_parser(item: Dict) -> List[Dict]:
        # Your custom parsing logic here
        return [processed_data]

datasets:
  my_dataset:
    type: file
    source: local
    path: "data/paths.json"
    parsing:
      - function: my_custom_parser
```