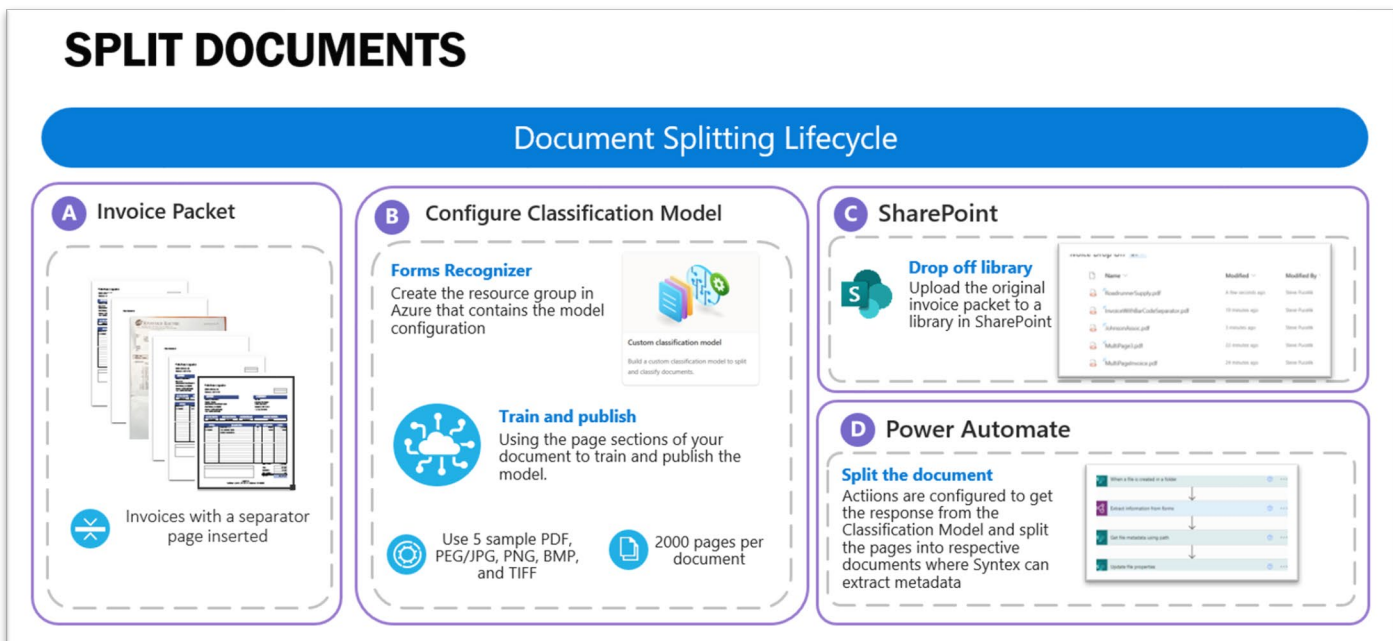# Split pages from documents

Often you will get a large document with several pages in it that are separated with a blank page, barcode or a change to the type of page that will distinguish a logical break in the document. This is common when you want to save storage space or put like documents together in a packet (mortgage applications, bundles of invoices, HR onboarding packets, health information etc.).

The newly released Custom Classification Model in Microsoft Forms Recognizer allows you to train a model on your documents to recognize the various portions of the document. This makes it easy for you to then break apart the document into logical smaller documents and then use Microsoft Syntex to classify and extract metadata.

## Process Overview

Here is a high-level overview of how to process the document and split the pages into smaller documents for processing.



## Structure of a document

To understand how to split pages from a document, you need to understand the structure of your document and what portions you need to train the Custom Classification Model with.

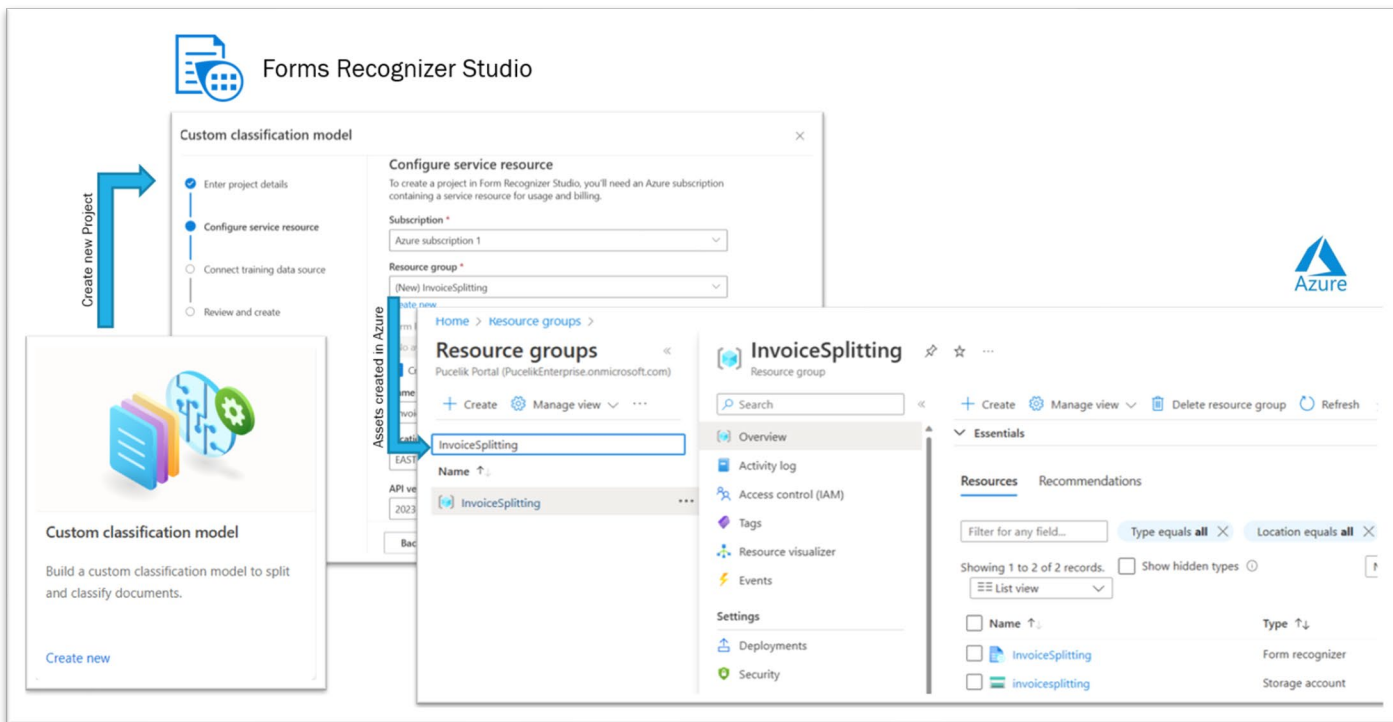Throughout this post I'll use a standard invoice as my example. The structure of the document will be:

- 2-page invoice (pages 1-2)
- Separator sheet (page 3)
- 1 page invoice (page 4)

As you train the Classification model, you will need at least 5 sample documents for each document section you want to split the pages on and create it's own document.

The end result is going to be 3 individual documents, a 2 page invoice, a separator sheet and a 1 page invoice. All of these will be created in a destination document library.
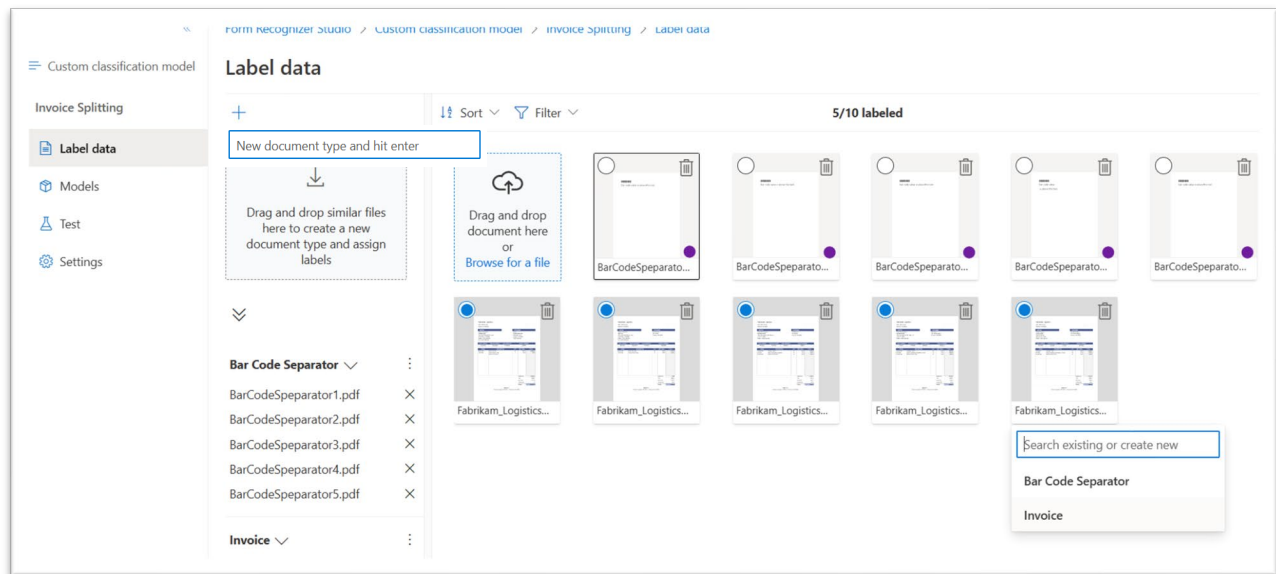
## Forms Recognizer Custom Classification Model

When creating a Custom Classification Model using Forms Recognizer Studio, the wizard will guide you through creating the appropriate configuration in your Azure portal. Once Forms Recognizer has finished setting up your project, you will see the components in Azure. This is where you'll get the endpoint and key needed when we create the Power Automate workflow.

## Labeling and training documents

Now that you have a project created, create the appropriate Document Types (Bar Code Separator and Invoice in the example below) and then upload 5 documents for each Document Type created.



Select the documents and associate them to the correct Document Type. Once you have all the documents labeled, you can train the model.

Once the model is trained, test it using a sample document with the page sections corresponding to the document types and you will see the results on the right. Also notice the Result tab, this is the JSON

output that we'll be using in Power Automate to split the documents.

## Power Automate Integration

Now that we have a trained model against the page sections in your document, we can use that in Power Automate against documents received in SharePoint. The complete solution can be downloaded from GitHub <INSERT LINK>. I'll walk through the critical actions and the configuration.



The highlighted areas indicate:

- This is the name of the model in Forms Recognizer Studio that was trained.
- The key from Azure can be found in the resource group that was created by Forms Recognizer.

- The initial request to Forms Recognizer will return a JSON file that contains the URL needed to retrieve the results once the document has been processed. This is in the header of the response named "Operation-Location".

- The delay activity gives Forms Recognizer time to process the document and return the JSON response. 30 seconds usually does it.
- Finally retrieve the JSON response containing the results we can use to split the document.
  - The yellow highlighted docType attribute indicates the document type Forms Recognizer classified it as based on how you trained the model.
  - The pageNumber represents the page in the document for the respective document type.



Logic configured in the Power Automate activities will evaluate each docType. When a docType changes from "Invoice" to "Bar Code Separator", we know to create a document that will contain a 2-page invoice starting on page 1 and ending on page 2. This process repeats until all docType attributes have been evaluated.

## Creating the new document

Once the Power Automate logic has determined that a new document needs to be created with the pages identified in the JSON file, I used the Split PDF action provided from Adobe to split the original file that was uploaded.
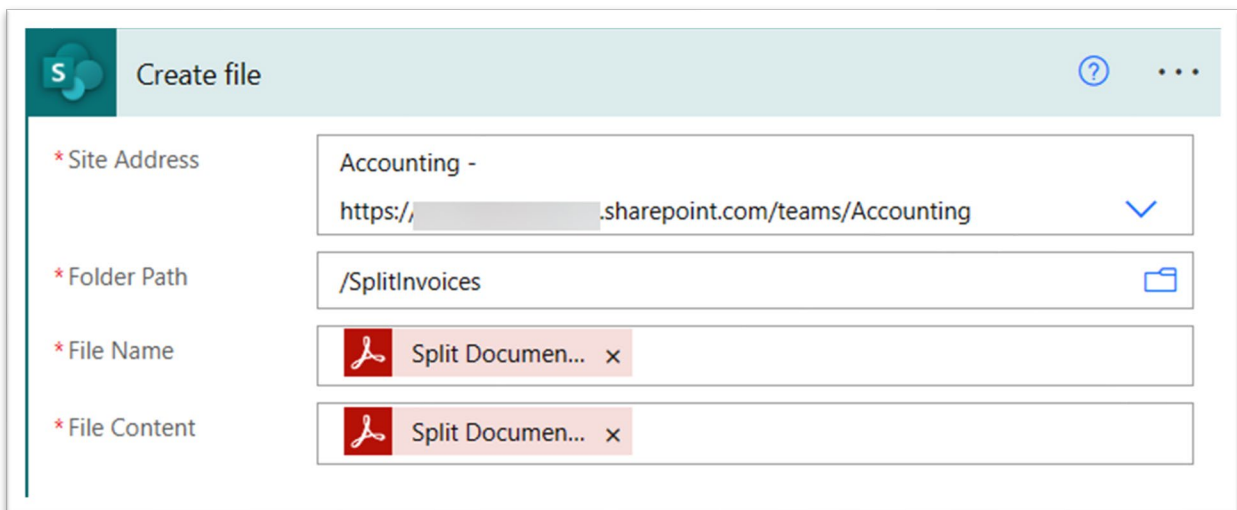


The resulting file is then saved to the final SharePoint library where a Microsoft Syntex model has been configured to classify and extract metadata from the file.

Environment Variables

- Drop off site
- Drop off library
- Destination site
- Destination library
- Model Name
- Ocp-Apim-Subscription-Key