

IBM Cloud Pak for Business Automation Demos and Labs 2022

Capture

IBM Automation Document Processing
V22.0.2

Lab Automation Document Processing

V 2.0

Clandis Baker
SWAT Business Automation Portfolio Specialist – Capture Products
bakercl@us.ibm.com

Krish Lakshminarayanan
Global Technical Program Leader for Capture / Intelligent Document Processing Global Sales (WW)
krishkrish@ibm.com

Ryan Sparks
Advisory Business Automation Tech Sales Leader – RPA/ADP
rmsparks@us.ibm.com

NOTICES

This information was developed for products and services offered in the USA.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing

IBM Corporation

North Castle Drive, MD-NC119

Armonk, NY 10504-1785

United States of America

The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law: INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some states do not allow disclaimer of express or implied warranties in certain transactions; therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this IBM product and use of those websites is at your own risk. IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements, or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.

TRADEMARKS

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is

available on the web at "Copyright and trademark information" at www.ibm.com/legal/copytrade.shtml.

Adobe, the Adobe logo, PostScript, and the PostScript logo are either registered trademarks or trademarks of Adobe Systems Incorporated in the United States, and/or other countries.

Cell Broadband Engine is a trademark of Sony Computer Entertainment, Inc. in the United States, other countries, or both and is used under license therefrom.

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

IT Infrastructure Library is a Registered Trade Mark of AXELOS Limited.

ITIL is a Registered Trade Mark of AXELOS Limited.

Java and all Java-based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

Linear Tape-Open, LTO, the LTO Logo, Ultrium, and the Ultrium logo are trademarks of HP, IBM Corp. and Quantum in the U.S. and other countries.

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

Microsoft, Windows, Windows NT, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

UNIX is a registered trademark of The Open Group in the United States and other countries.

© Copyright International Business Machines Corporation 2020.

This document may not be reproduced in whole or in part without the prior written permission of IBM.

US Government Users Restricted Rights - Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

Table of Contents

1.	Overview.....	6
1.1	Getting HELP during the lab	6
1.2	Icons	6
1.3	Abstract.....	6
1.4	Introduction	7
2	Getting started	8
2.1	IBM TechZone – Overview	8
2.1.1	Reserve Environment	9
2.2	Open your IBM Cloud Environment	12
3	Lab Overview.....	14
3.1	How does ADP work?	14
4	Create Document Processing Project.....	16
4.1	Reviewing the interface.....	21
4.1.1	Build Tab	21
4.1.2	Enrich Tab	22
4.1.3	Configure Tab.....	23
5	Configure a Wage and Tax document type.....	26
5.1	Create Wage and Tax document type.	26
5.2	Create Field	28
5.3	Create the Employee Name Address field.....	33
5.4	Create Employee Social Security Number Field	34
6	Document Types and Samples Overview	37
6.1	Categorize documents.	37
7	Train classification	45
7.1	How do I improve my results?	50
7.1.1	Option 1 – Add more samples.	50
7.1.2	Option 2 – Review all uploaded samples.	51
8	Data extraction.....	52
8.1	Correcting extracted values	56
8.2	Train extraction model.....	61
9	Data standardization	62
10	Version and deploy your project.....	64
11	Application designer	67
11.1	Create your Runtime Application.	67
11.2	Upload documents for processing.....	74
11.3	Correct any classification errors.....	77
11.4	Correct extraction issues.....	79
12	Optional Export/Import Project.	85
Appendix A - Troubleshooting	87	
TechZone Pending Status taking Long Time.....	87	
Service Error.....	87	
Application Blank	88	

Connection issue with Workstation to Cloud.....	89
Opening an Incognito Window	89
Popup Blocked when trying to Preview Application.....	90
Appendix B - BAW & ADP Integration Sample	91
Appendix C - Badge Information.	92

1. Overview

1.1 Getting HELP during the lab

- For internal IBM, another good resource is the Archive slack channel for questions:
#cp4ba-adp-lab or <https://ibm-cloud.slack.com/archives/C01LVVBMWPN>
- For external participants besides the Slack channel, use the Webex chat if you are in a webex event or just speak up.
- For others, email bakercl@us.ibm.com. This method will be slower and will be best effort. It may require jumping on a Webex meeting to provide help.
- Getting help after lab reach out to the following:
 - bakercl@us.ibm.com
 - krishkirsh@us.ibm.com
 - rmsparks@us.ibm.com

1.2 Icons

The following symbols appear in this document at places where additional guidance is available.

Icon	Purpose	Explanation
	Important!	This symbol calls attention to a particular step or command. For example, it might alert you to type a command carefully because it is case sensitive.
	Information	This symbol indicates information that might not be necessary to complete a step but is helpful or good to know.
	Trouble-shooting	This symbol indicates that you can fix a specific problem by completing the associated troubleshooting information.

•

1.3 Abstract

Set up a capture solution in minutes. Introduce technical sellers to IBM Automation Document Processing. In this session, students will configure their own capture project. They will learn how to use machine learning classification for their sample documents, define fields for extraction, create validation rules, and use deep learning* (subject to environment configuration) to automate data extraction.

1.4 Introduction

Welcome to the Automation Document Processing lab. This lab will introduce you to Document Processing and provide you with an understanding how you can configure it for your customer opportunities.

Automation Document Processing provides a tailored solution that reads your documents (in English, French, Spanish, German, Dutch, Portuguese), extracts data, and refines and stores the data for use.

With the right business knowledge, you can design deep learning models without being a data scientist. The Document Processing Designer includes pre-trained deep learning models that you can use as a base for your own model. The pre-trained document types include bills of lading, invoices, and utility bills.

You can extract text, check boxes, forms, tables, barcodes, signature detection and even free text. With no or low code options, you can create an application that processes documents, extracts data, flags issues, and stores your documents and data. And the data enrichment capabilities ensure that the extracted data is standardized and ready for use in downstream integrations.

This lab will not cover all the available functionality available due to time constraints. Additional labs will be created in the next few months to add to your knowledge and understanding of Document Processing.

2 Getting started

Download the sample documents in the zip file. We will be using these sample documents during the labs You can find them here:

<https://github.com/IBM/cp4ba-labs/tree/main/22.0.2/Document%20Processing>

File	Action	Last Modified
Group 1 - Design Docs for Tax Lab.zip	Add files via upload	1 hour ago
Group 2 - Classification Results Increase Set.zip	Add files via upload	1 hour ago
Group 3 - Runtime demo Set.zip	Add files via upload	1 hour ago
Readme.md	Update Readme.md	5 days ago
[In Process]Lab Guide - Automation Document Proce...	Add files via upload	4 days ago

_1. Click on “Group1 – Design Docs for Tax Lab.zip”.

_2. Then Click on Download

Code Issues Pull requests Actions Projects Wiki Security Insights

1.23 MB

View raw

Download

_3. Repeat above steps “Group 2 – Classification Results Increase Set.zip” and “Group 3 – Runtime Set.zip”

_4. Unzip the files and keep them in their designated folder.

You will notice the images are in various unique folders that will be referenced specifically in the different labs later. Please keep them in their proper folders.

2.1 IBM TechZone – Overview

What is IBM TechZone?

IBM Technology Zone (techzone.ibm.com) enables IBM teams and IBM Business Partners to provision technical “Show Me” live environments, Proof-of-Technologies, prototypes, and Minimum Viable Prototypes, which can be customized, shared with peers and clients to experience IBM Technology.

Learn more: <https://techzone.ibm.com/collection/onboarding#tab-1>

The TechZone leverages DAFFY. DAFFY is Deployment Automation Framework For You. The DAFFY installer tool has been renamed to Pak Installer. This tool will do all the heavy lifting of the OpenShift and IBM Cloud Pak installs. The National Market Top Team created Pak Installer to assist the technical sales teams with the progression of IBM Cloud Pak opportunities.

The goal is to provide the technical sales with a set of (easy to use) scripts that will aid in the installation of OpenShift and the IBM Cloud Pak's. For more information on DAFFY/Pak Installer please look at: <https://ibm.github.io/daffy/Deploying-OCP/TechZoneTilesBeta/>

Also please rate the tile on the TechZone website. This will help the development team understand its value.

2.1.1 Reserve Environment

- _1. Navigate to <https://techzone.ibm.com/collection/63457fcba311ed0018ca2442>



Note: Don't grab just any CP4BA in TechZone. This environment is built with Daffy authored by Kyle Dawson with the latest releases. This environment can also be used at a customer site with same tool and framework of Daffy.

- _2. Click Cloud Pak for Business Automation tab and scroll down to the “Cloud Pak for Business Automation 22.0.2 – VMWare tile.
- _3. Click on Reserve

- _4. On Create a reservation screen **select option** for when to start

- _5. Create a Reservation

Based on the reservation type you are making, provide the required information

Customer Demo : Need a short customer-facing demonstration

Practice/Self-Education: Need to gain experience

Standard proof of concept; Need an environment for a standard product use case.

Custom Proof of concept: Need a complex, customized environment.

Testing: Need to test a specific function, configuration, or customization.

- _6. For this lab **Select Testing** will give you 3 days plus the option to extend it for another week. Otherwise, you will need a legitimate opportunity to leverage another reservation type.

- _7. For Preferred Geography (required) select your preferred data center location

- _8. For VPN Access **choose Enable**. You will be using a VPN to connect from desktop to the TechZone tile



Make sure to pick enable otherwise you'll have to start all over with deployment.

- _9. In Cloud Pak for Business Automation Version Pick 22.0.2. (if not already chosen for you)

- _10. In Cloud Pak for business Automation IFix pick IF002 (if not already chosen for you)

- _11. For Starter Service **choose docprocessing**

12. Click Submit

You will receive emails on the status and the steps that have completed and please allow up to 1 hour for the start-up services to fully complete. Your emails will look like the following:

IBM Technology Zone

Status Update:
Cluster Running, Now Starting Cloud Pak

Step	Status	Step ETA
Cluster Installed	Complete	90 Min
Cloud Pak Installed	In Progress	60 Min
Cloud Pak Services Installed	Pending	60 Min
Cloud Pak Services Running	Pending	120 Min

To connect to your Cluster, logon to the Portal to get details about your environment.

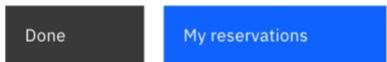
PakInstaller Portal - <http://connection-info.apps.mycluster.ibm.com:32080>
User ID - pakinstaller
Password - mypassword

IBM Technology Zone - Brought to you by PakInstaller

Once “Cloud Pak Services Running” is Complete upon receiving this email then your environment is ready. Once the start-up process is complete you can click on the links identified in the email. However, it is recommended that you review your reservation information from the IBM Technology Zone – My reservation site.



If after receiving email and a few hours have passed and your environment is not up, check [Appendix A – Troubleshooting](#) for possible fix.

13. Click My reservations**14. Once you get the email from the IBM Technology Zone site, you can access your environment reservation(s) by clicking on the **My library** then **My Reservations****

IBM Technology Zone | My library ^ | Help

Welcome,

Get started

Updated Nov 13, 2022

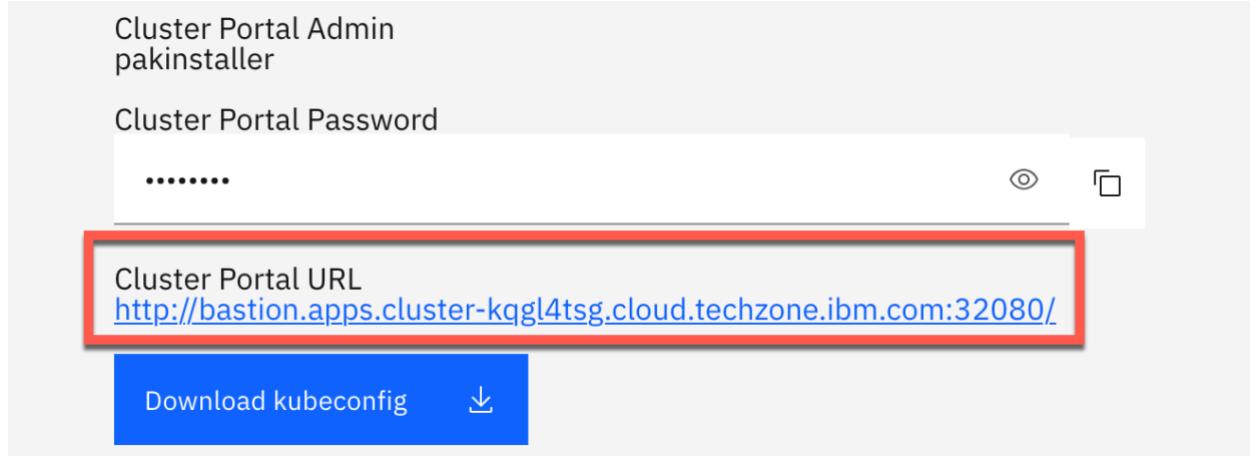
IBM Technology

- My created collections
- My bookmarked searches
- My favorites
- My reservations**
- My workshops

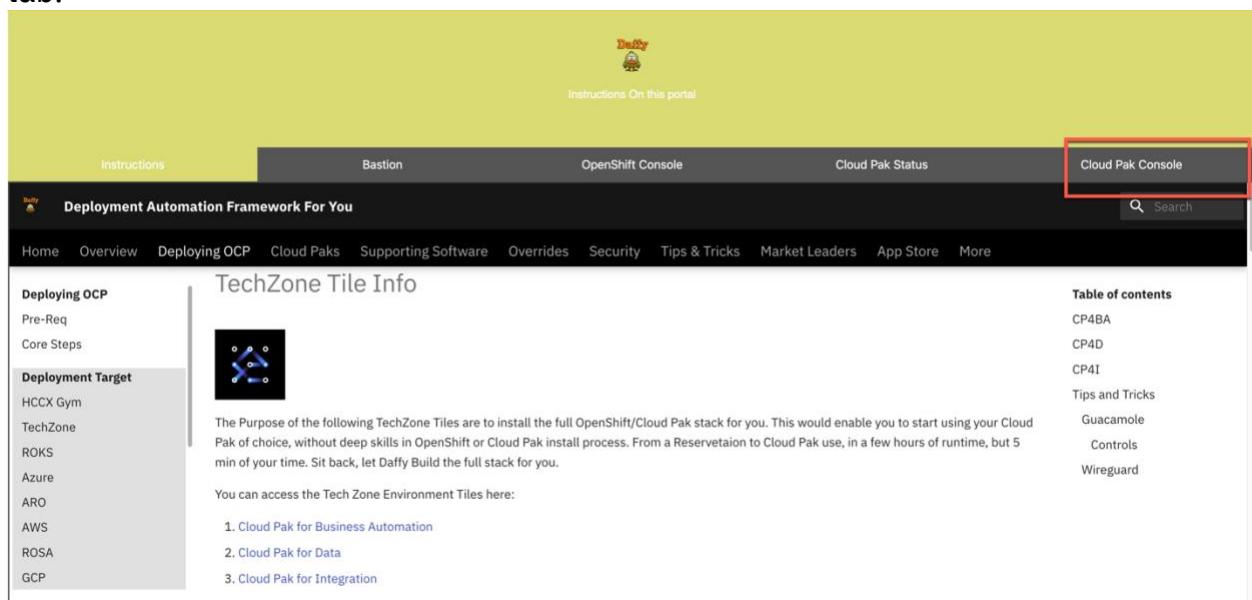
You can also access directly using the link below
<https://techzone.ibm.com/my/reservations>

2.2 Open your IBM Cloud Environment

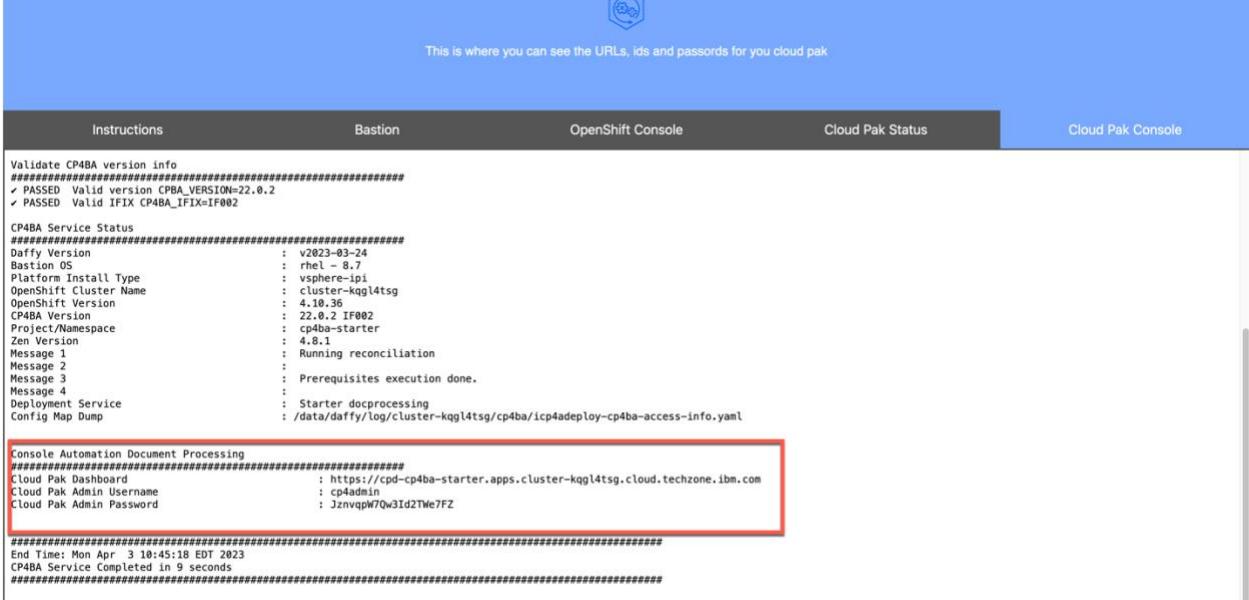
- _1. Back on your reservation screen under the Environment section **Right Click** on the link under **Cluster Portal URL** and **Select open in a private window**



- _2. Login with <**Cluster Portal Admin /Cloud Pak Admin Password**> provided in Environment section above. In example above pakinstaller/<your unique password>.
- _3. You will be presented with the Daffy portal screen. Click on **Cloud Pak Console** tab.



- _4. Under the Cloud Pak Console tab scroll down to Console Automation Document Processing



This screenshot shows the 'Cloud Pak Console' tab selected in a browser. The page displays various status metrics and logs. A red box highlights the 'Console Automation Document Processing' section, which contains the URL for the Cloud Pak Dashboard: <https://cpd-cp4ba-starter.apps.cluster-kqgl4tsg.cloud.techzone.ibm.com>. Other visible log entries include Bastion OS details, OpenShift cluster information, and deployment service logs.

```

This is where you can see the URLs, ids and passwords for your cloud pak

Instructions Bastion OpenShift Console Cloud Pak Status Cloud Pak Console

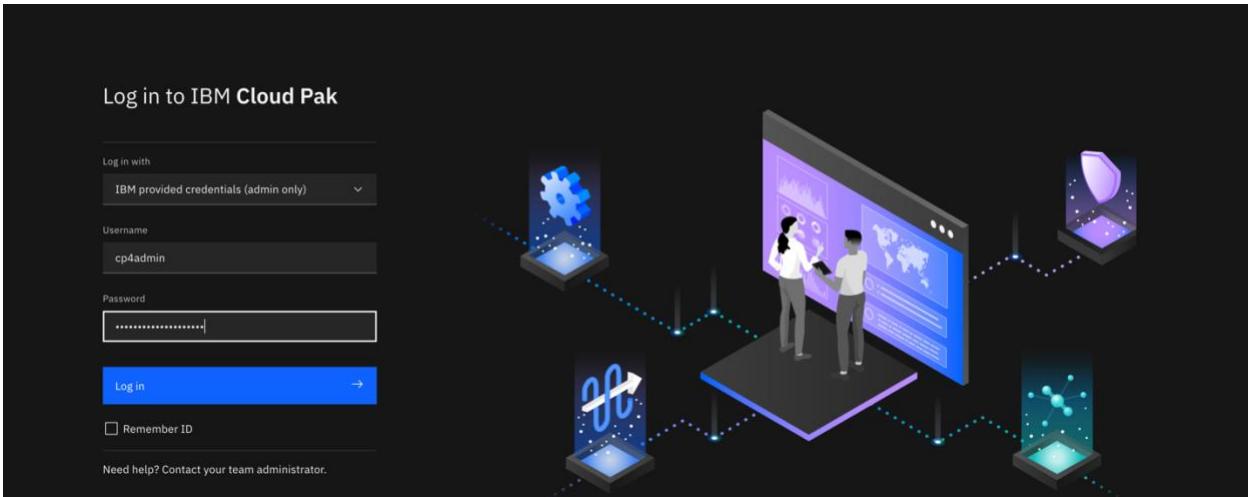
Validate CP4BA version info
#####
# PASSED Valid version CP4BA_VERSION=22.0.2
# PASSED Valid IFIX CP4BA_IFIX=IF002

CP4BA Service Status
#####
Daffy Version : v2023-03-24
Bastion OS : rhel - 8.7
Platform Install Type : vsphere-1pi
OpenShift Cluster Name : cluster-kqgl4tsg
OpenShift Version : 4.13.36
CP4BA Version : 22.0.2 IF002
Project/Namespace : cp4ba-starter
Zen Version : 4.8.1
Message 1 : Running reconciliation
Message 2 :
Message 3 : Prerequisites execution done.
Message 4 :
Deployment Service : Starter docprocessing
Config Map Dump : /data/daffy/log/cluster-kqgl4tsg/cp4ba/icp4adeploy-cp4ba-access-info.yaml

Console Automation Document Processing
#####
Cloud Pak Dashboard : https://cpd-cp4ba-starter.apps.cluster-kqgl4tsg.cloud.techzone.ibm.com
Cloud Pak Admin Username : cp4admin
Cloud Pak Admin Password : Jznvpgk7Qw3Id2TWeFZ
#####

End Time: Mon Apr 3 10:45:18 EDT 2023
CP4BA Service Completed in 9 seconds
#####
  
```

Copy the link to your favorite private browser or open in new tab. Accept any potential security risks. You will then be prompted with Log in screen to the IBM Cloud Pak.



- _5. Enter the <**Cloud Pak Admin Username**> and <**Cloud Pak Admin Password**> from Cloud Pak Console tab. **Click** on blue **Login** button.
- _6. You will be presented with the “Welcome! Let’s get started” screen. **Click** Maybe later button.



Note you will see this screen several times throughout the lab. You can always select Maybe later while doing this lab.

3 Lab Overview

The lab will focus on the design time tasks for Automation Document Processing (ADP). Despite the push for the digitization of content for many years, there are still a lot of paper documents that require workers to read and interpret the information – whether it is structured data, such as tax forms, or semi-structured data, such as invoices, utility bills, and so on. This lab describes how to set up an automate document processing pipeline using ADP.

3.1 How does ADP work?

Document Processing Designer

You use the Designer interface to create a set of document types and related fields that comprise your Document Processing project. Document Processing Designer combines an intuitive interface with a set of AI and deep learning tools that identify and learn the document types that matter to your organization. For each document type, you designate which pieces of information to extract as data for that document to be used by downstream applications. You can also apply tools to clean up and standardize the data as it is extracted.

Deployment tools

After you build the Document Processing project in the Designer, you deploy the project to make it available for building your document processing application. The deployment process is also used to configure the repository to receive the processed documents from your end-user application by making the capabilities and artifacts available for integration into an application and into the destination repository.

Application templates and toolkits

You use the no- or low-code application building capabilities of Application Designer, customized templates and toolkits, and the AI model of your Document Processing project to create a document processing end-user application. This application recognizes your documents, extracts your relevant data, and presents issues to fix before sending the documents to storage and using the data in other systems.

Document processing application and document management

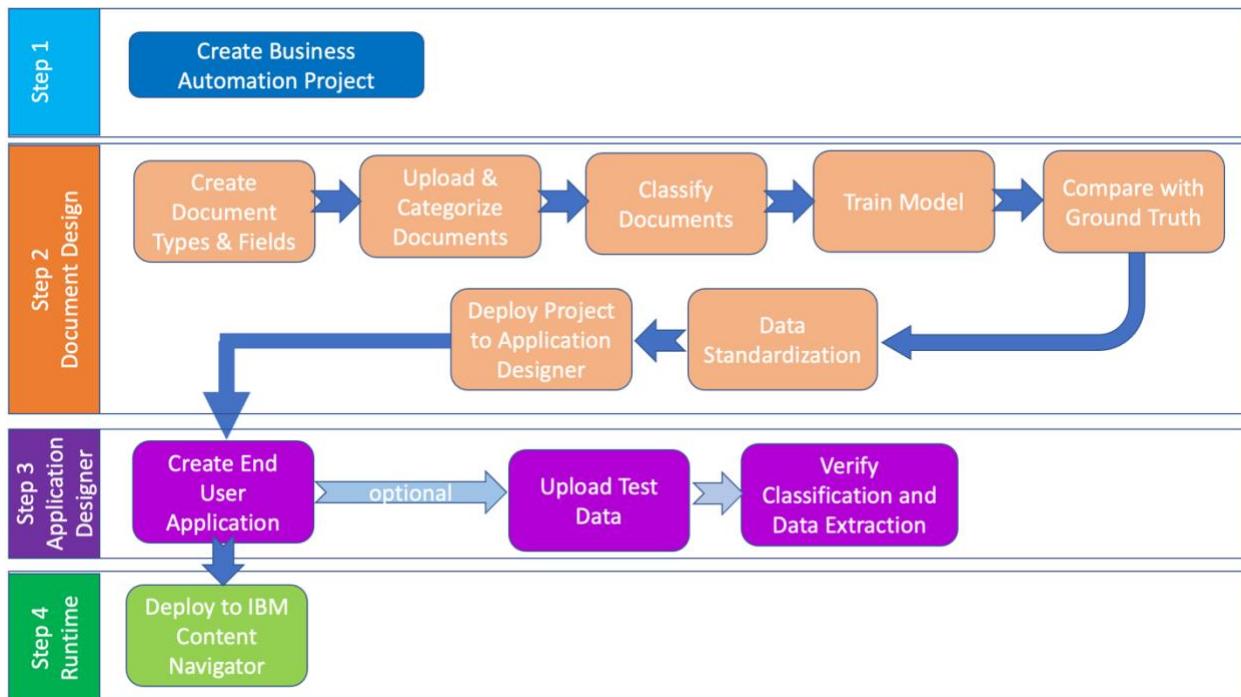
The application that you build uses the AI and deep learning to automatically detect, extract, and standardize the data in all your documents. Any anomalies are flagged according to your customized model and the priority that you set so that your document processing user can correct issues before the documents are finalized.

When you deploy your document processing application, you connect it to a content repository that manages the document types and the extracted data for each document. The solution is fully integrated with IBM FileNet® Content Manager, simplifying document and data storage by applying your existing filing architecture and business rules to each processed document. The content and metadata are automatically saved in FileNet within the appropriate document class.

End result

Your document types are stored in the content repository, with appropriate retention and access controls. An associated JSON file reflects all the extracted data for the document. Properties are set on the document with the data definition-controlled values. Your extracted data is cleaned, standardized, and ready for use in other applications.

The following diagram shows the tasks required to configure and deploy a new ADP project.



Step 1 – Create an ADP Business Automation Project

Each document processing project requires a separate repository in your Git organization. Coordinate with your Git administrator to create the repository for your project.

Step 2 – Document Design

This step shows the high-level tasks that will be needed to complete to train the system to recognize document types, successfully extract fields and tables, configure the fields in FileNet and finally deploying your content project to the application designer so you can configure the end-user interfaces.

Step 3 – Application Designer

The application designer is where you would configure end-user interfaces such as the classification and verification screens. The lab will not go in a lot of details on how to configure the interfaces. It will instead show you how to create an application, and test processing a batch of documents through the system. To get more information on creating/using the Business Automation Application (BAA) look at the SWAT Jam Lab for BAA.

Step 4 – Runtime

End-users would be using the runtime IBM Content Navigator interface to process documents or batches, classify document and verify extracted field data in the verification screen.

4 Create Document Processing Project

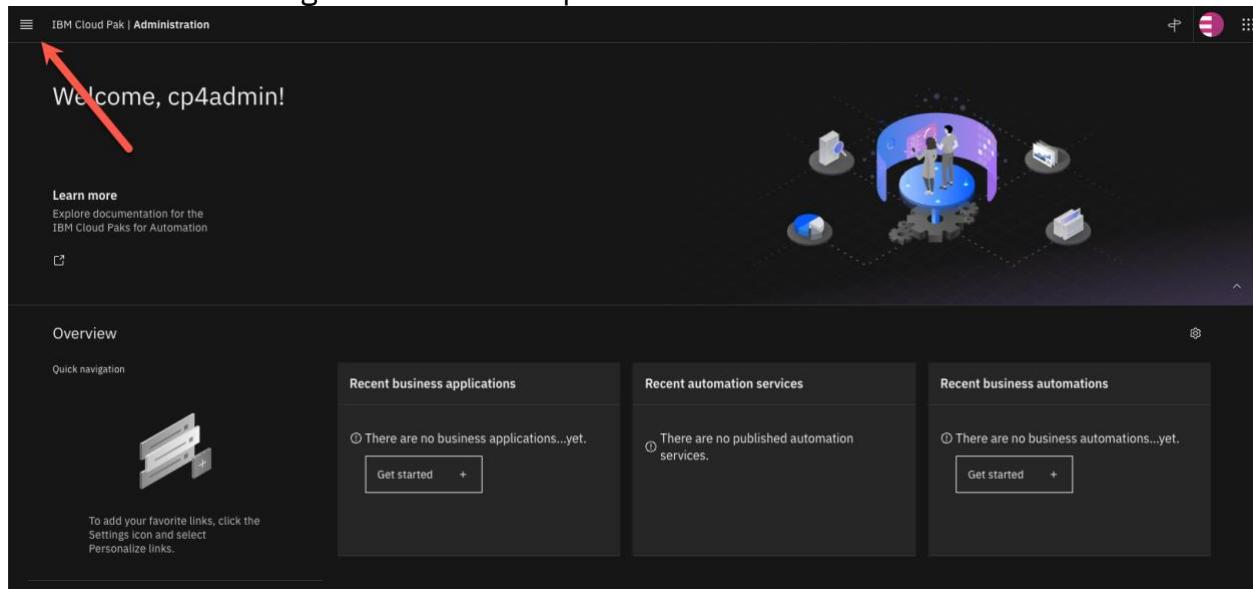
Step 1

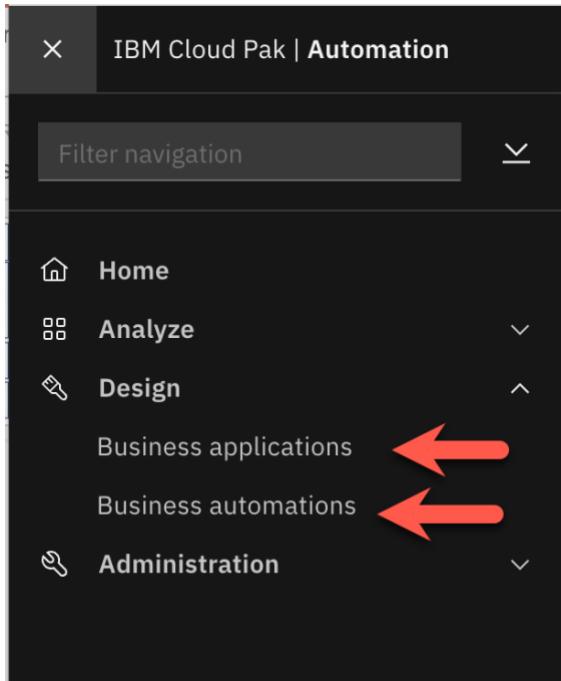
Create Business Automation Project

Cloud Pak for Business Automation Studio is the single authoring and development environment for the IBM Cloud Pak for Business Automation platform that accelerates digital transformation. Business Automation Studio provides an entry point to various designers to help you reach your goals.

There are two distinct parts to the Business Automation Studio configuration.

- _1. Click on the hamburger menu at the top left next to IBM Automation.





Business automations provides the Document Processing configuration of the document classes, and the **Business applications** provides the user interfaces.

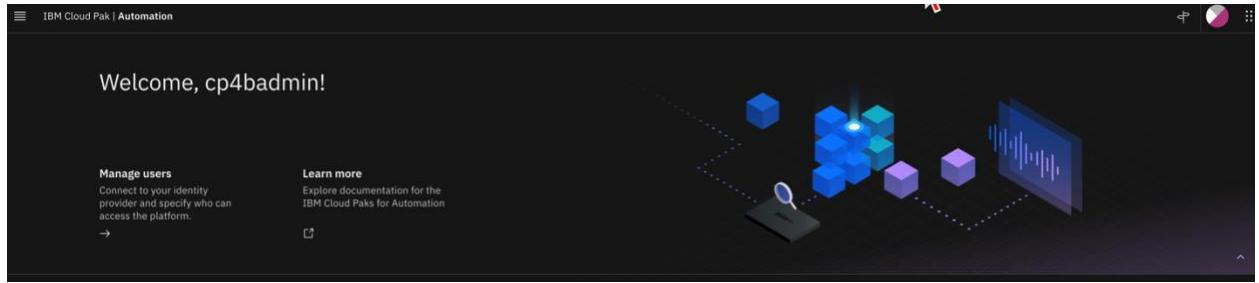
Within the Business Automations you can create or reuse automations. An automation is a collection of artifacts that fulfills a business purpose. You can publish some automation artifacts as automation services that you can call and reuse in a consistent way. Also in Business Automation, you use the **Document Designer** interface within Automations to create a set of document types and related fields that comprise your Document Processing project.

The Document Processing Designer combines an intuitive interface with a set of AI and deep learning tools that identify and learn the document types that matter to an organization. For each document type, you designate which pieces of information to extract as data for that document to be used by downstream applications. You can also apply tools to clean up and standardize the data as it is extracted.

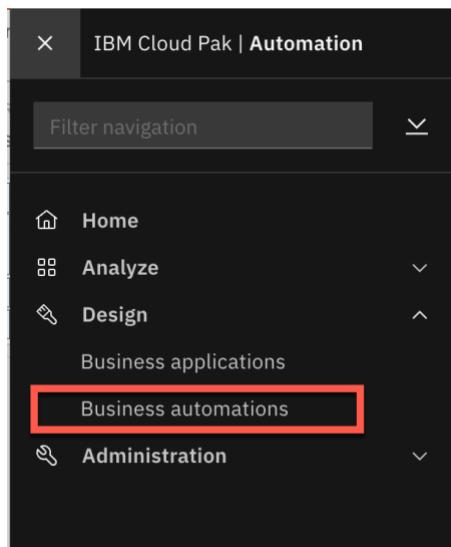
Within **Business Applications** you can quickly create user interfaces that integrate tasks, data, and automations. You can start with a template to ensure consistency. You can also use toolkits to share artifacts from existing applications.

We will start with the Business Automations.

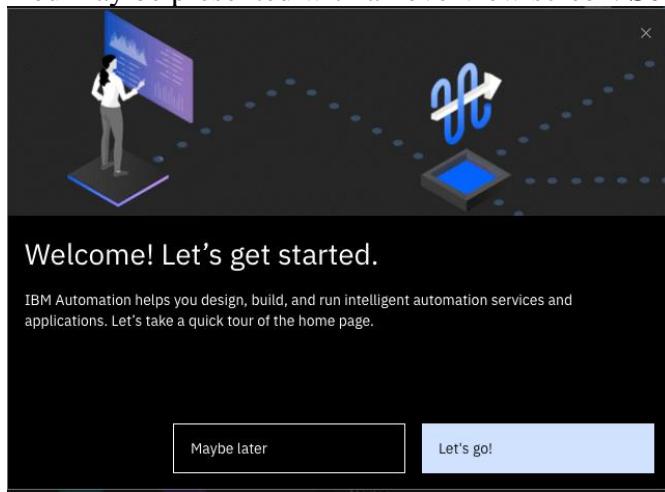
Once logged in to the IBM Automation Server, you should see the Welcome screen.



_2. Click on Drop down arrow next to Design then Select Business Automations.



You may be presented with an overview screen. **Select Maybe Later.**



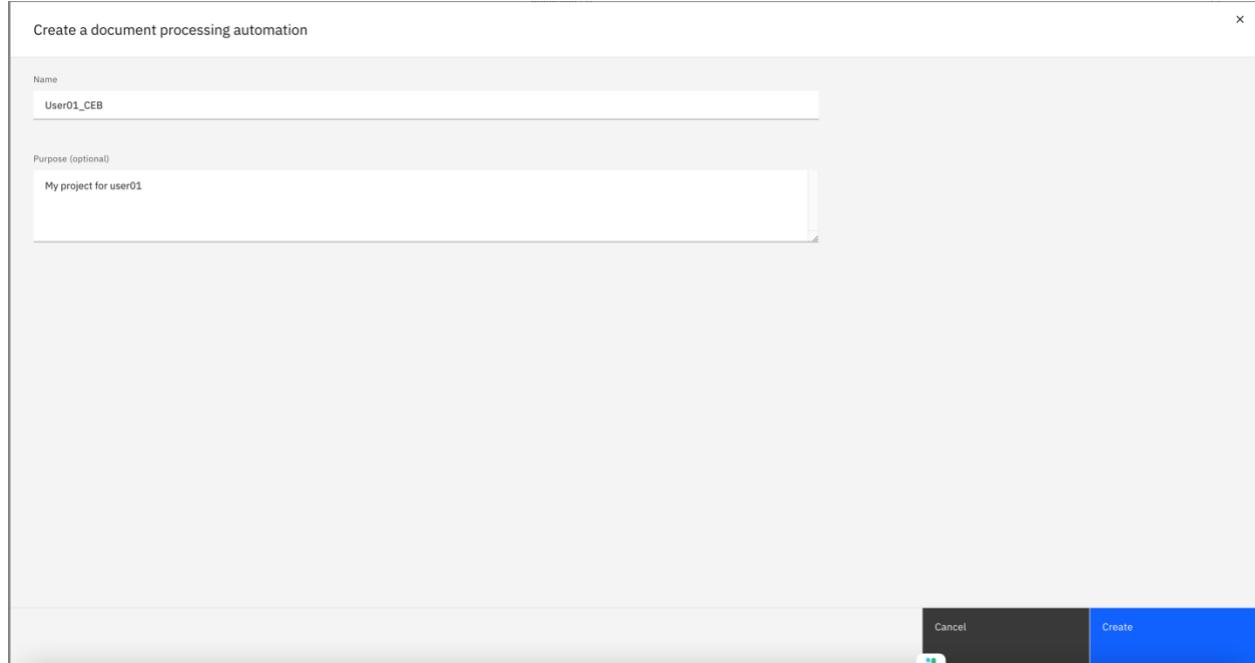
Then following screen appears.

The screenshot shows the 'Business automations' page. At the top, there's a navigation bar with the text 'IBM Cloud Pak | Automation'. Below the header, the title 'Business automations' is displayed. A brief description follows: 'Create or reuse automations. An automation is a collection of artifacts that fulfills a business purpose. You can publish some automation artifacts as automation services that you can call and reuse in a consistent way.' A 'Learn more' link is present. Below the description are two main buttons: 'Create' (highlighted with a blue background) and 'Import'. Underneath these buttons is a list of categories: 'Published automation services' (with a right-pointing arrow), 'Decision' (with a right-pointing arrow), 'Document processing' (with a right-pointing arrow), 'Workflow' (with a right-pointing arrow), and 'External' (with a right-pointing arrow).

7. Click on the **Create** twisty and select **Document processing automations**.

This screenshot shows the 'Create' dropdown menu open. A red arrow points from the previous screenshot's 'Create' button down to this menu. The menu items listed are 'Decision automations', 'Document processing automations' (which is highlighted with a red box), 'Workflow', and 'External'. Below the dropdown, the main 'Business automations' page is visible again, showing the same structure as the first screenshot but with the 'Create' menu open.

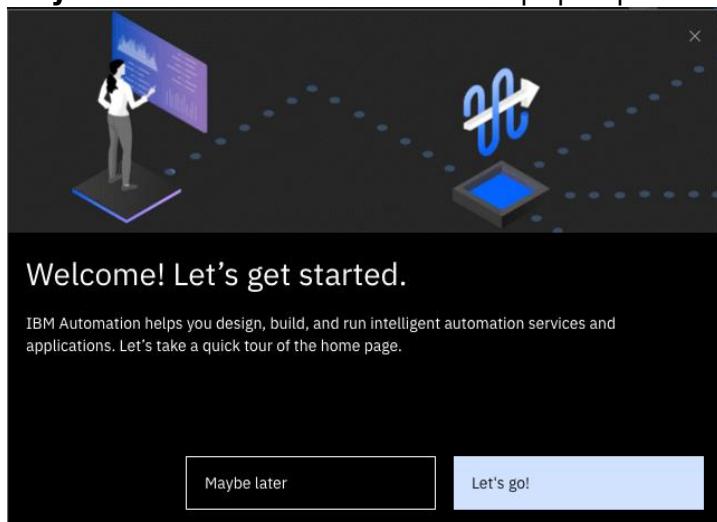
- _8. In the Create a document processing automation window **enter a name** for the project. Optional, enter a purpose.



- _9. Click on **Create** in the lower right-hand corner.



You may see the *Welcome Let's get started* throughout the lab simply **click Maybe later** whenever this window pops up.



4.1 Reviewing the interface.

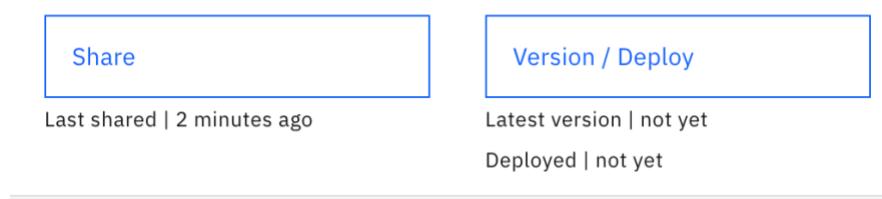
The screenshot shows the IBM Cloud Pak | Administration interface for a project named "Clandis Baker Project". The "Build" tab is active. The main content area displays five configuration sections:

- Document types and samples:** Shows 3 types and 29 samples on average. Status: Ready.
- Classification model:** Shows 3 types trained with 100% accuracy. Status: Ready.
- Extraction model:** Shows 3 types trained with 97% accuracy. Status: Ready.
- Data standardization:** Shows 0 types mapped. Status: Not ready.
- Document retention:** Shows 3 types reviewed. Status: Ready.

On the top right, there are "Share" and "Version / Deploy" buttons. The "Share" button indicates "Last shared | a few seconds ago". The "Version / Deploy" button indicates "Latest version | not yet Deployed | not yet".

Upon opening the project, there are three major sections: **Build tab, Enrich tab, and Configure tab.**

On the top right, you find the *Share* and *Version/ Deploy* buttons.



The *Share* button is used to save your configuration to your GitHub repository.

The *Version / Deploy* button is used to create a snapshot, or version of your configuration. Like the *Share* button, the *Version* button will save your configuration, but will also create a version of it while retaining your previous version.

Once you have created a version of your configuration, you can also use this button to *Deploy* your version to the Business Applications area of ADP. You need to do this before you can go into the Business Application tile and configure your user interfaces.

4.1.1 Build Tab

This is what you will be spending most of your time on. The BUILD tab shows the guided configuration for building a Document Processing project. It shows the five steps required.

Document types and samples: Here you will define the document types that can be recognized by this automation and upload sample documents for training. By default, any project will be pre-populated with three pre-trained document types (Bill of Lading, Invoice, and Utility Bill).

Classification model: classification: Here we will teach the system how to recognize the different document types.

Extraction model: Here we will teach the system how to extract information for each document type based on the classification.

Data Standardization: This allows further refinement of the extracted information. For example, we want to standardize all dates to be formatted as YYYY/MM/DD. Having a standardized data format will help with any subsequent automation process.

Document retention: This allows us to define how long we want our documents to be kept in the system. Documents that have exceeded the retention period will be automatically expunged. This could be important for regulatory compliance or for managing the overall storage size.

4.1.2 Enrich Tab

_1. Click on the ENRICH tab.

Enrich provides a quick way to define your document types and the fields you wish to extract. In this section, we can define additional enrich rules. An example of an enrich rule is to specify the expected format for an invoice number (all numerical) or a driver's license. The more we can tell document processing about how different data will be formatted, the higher the chance it will recognize the information.

_2. Click on FIELD TYPES AND ENRICHMENTS to begin. In this tile, you will see some of the pre-configured fields in the SYSTEM LIBRARY (sys). Customers can use

these fields in their document type field definitions as needed.

Field type	Value type
Address block	String
Address information	Composite
Addressee	String
Boolean	Boolean
Building number	String
City	String
Country	String
Country code	String
Country name	String
Currency	Composite
CurrencyCode Object Type	String
Date	Date
Date Range	Composite
Decimal	Decimal

- _3. Click on <**your project name**> in the bread crumb trail at the top to go back to the Build tab.

IBM Cloud Pak | Administration
Business automations / Clandis Baker Project / Field types and enrichments

4.1.3 Configure Tab

- _4. Click on **Configure** Tab

This is where we can configure other operational aspects of the project. The export project creates a .zip file that contains the document types, field types and enrichments, which you can use to start training with new sample files. You can also decide to include the training model and the sample training files in your export if you want to move your entire project to a new instance of Document Processing for example. To import a project, select the .zip file to import. When you import a .zip file you have two options: overwrite the existing project or merge the existing project. If you merge the existing project, document types, field types, enrichments, and sample training files are imported unless there is a conflict. Models are not imported.

The screenshot shows the 'Configure' tab selected in the navigation bar. Under 'Import / Export ontology', there are three sections: 'Language settings' (selected), 'Git server configuration', and 'Export project'. The 'Export project' section contains a button labeled 'Export project'. Below it is the 'Import project' section, which includes a note: 'You cannot import a project if the current project has been deployed.' and a button labeled 'Import project'.

In Extraction language, select which languages are used in the documents that you plan to process. You can choose English, Dutch, French, German, Brazilian Portuguese, or Spanish. Make sure to choose only the language or languages that are likely to be used in your document sets. Choosing more than one language can affect the accuracy of your document processing model.

In Display name language, select the language that you use to enter display names for fields and document types. These are the names that are displayed in the Designer and in the applications.

The display name language is also used in the Content Engine as the localized string locale setting for document classes and properties. Document Processing project deployment supports only one language per project. If your organization has multiple projects with different language settings, these projects cannot be deployed to the same Content Engine server if they share common properties. For example, when you define data definitions during data standardization, you cannot map a field to an existing data definition that was created in a different language.

The screenshot shows the 'Configure' tab selected. In the 'Language settings' section, under 'Extraction language', there is a note: 'Choose the language you want to use for extraction. You can choose more than one language, but this will impact accuracy. [Learn more](#)'. A checkbox for 'English' is checked. Under 'Other languages', there are checkboxes for Dutch, French, German, Portuguese (Brazil), Spanish, and German. In the 'Display name language' section, there is a note: 'Choose which language you will use to create display names throughout the project. If you are naming items in this project in a language other than English, set the display name language to the language you will be working in.' A dropdown menu for 'Project locale' shows 'English (en) (default)'. At the bottom right are 'Cancel' and 'Save' buttons.

The Git server configuration is where you create a connection to the Git server for the first project that you create in Document Processing Designer. This setting applies to all

subsequent projects that you create.

IBM Cloud Pak | Administration

Business automations / Clandis Baker Project

Build Enrich Configure

Import / Export ontology

Language settings

Git server configuration

In order to share, version and deploy, you need to establish a connection to your organization's Git server.

Git vendor: Gitea

Git server organization URL: `https://cp4adeploy-gitea-svc:3000/content-designer`

Git server REST API URL: `https://cp4adeploy-gitea-svc:3000/api/v1`

Username: git

Type of credentials: API key Password

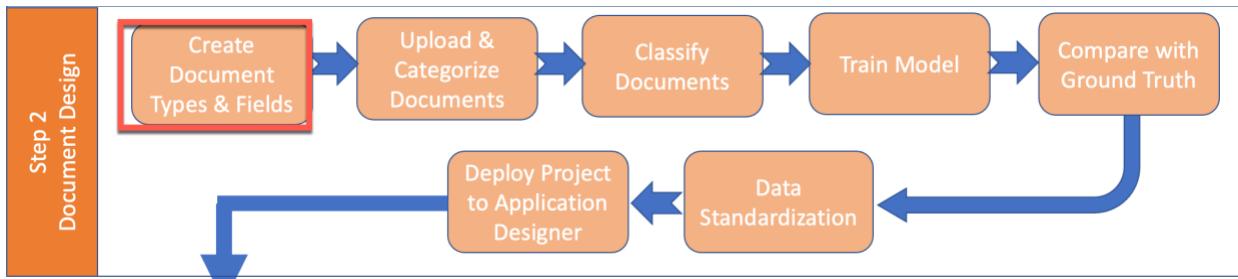
Credentials: Enter a password or API key

Test Save

Share Last shared 1 day ago

Version / Deploy Latest version v2 1 day ago
Deployed v2 1 day ago

5 Configure a Wage and Tax document type.



Before we use the guided configuration, you will configure some additional document types and fields used to extract data prior to uploading sample documents.

To do this lab, we will use the ENRICH tab to add fields to a newly created Wage and Tax document type.

5.1 Create Wage and Tax document type.

_1. Click on **Build** tab to return to the start page.

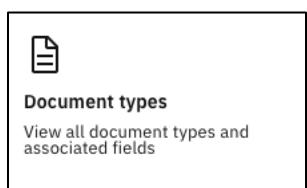
The screenshot shows the IBM Cloud Pak | Administration interface. At the top, there is a navigation bar with 'IBM Cloud Pak | Administration' and various icons. Below it, a breadcrumb trail shows 'Business automations / Clandis Baker Project'. On the right side, there are 'Share' and 'Version / Deploy' buttons. The main content area has three tabs: 'Build' (which is highlighted with a red box), 'Enrich', and 'Configure'. The 'Build' tab displays several status cards:

- Document types and samples**: Shows 3 types and 29 samples on average.
- Classification model**: Shows 3 types trained with 100% accuracy.
- Extraction model**: Shows 3 types trained with 97% accuracy.
- Data standardization**: Shows 'Not ready'.
- Document retention**: Shows 3 types reviewed.

At the bottom right of the interface, there are buttons for 'Last shared | 4 minutes ago' and 'Latest version | not yet Deployed | not yet'.

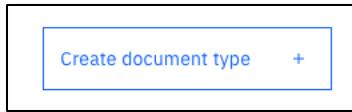
_2. Click on the **ENRICH** tab

_3. Click on **DOCUMENT TYPES**



We will now create a document type for Wage and Tax documents and fields to extract data from them.

_4. Click on the **CREATE DOCUMENT TYPE** button in the top right corner.



- _5. The Add document type window pops up. **Enter Wage and Tax** for the display name. There is no need to enter a symbolic name ADP will use the display name as a base. There's no need to add description in this lab unless you want to.

Add document type X

Display name 12/50
Wage and Tax
This is the name that will show up for you in the system. You can use characters from any language.

Symbolic name 10/50
WageandTax
This name will be used to identify the document type in the code.

Description (optional) 0/512
Enter a description for this document type

Fixed-format document type
Fixed-format documents have a fixed structure that remains the same for every document. Fixed-format documents types do not require as many sample documents to be trained in the extraction model.
 This document type has a fixed format

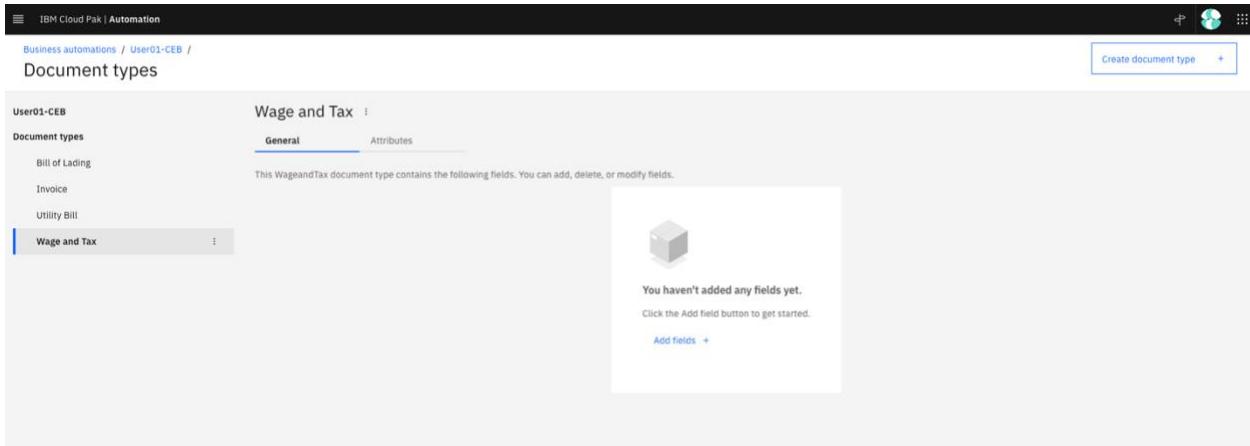
Cancel **Add**



Note: Notice the option for “Fixed-format document type”. If your form is static in nature or has a fixed structure that does not change, select this option so you will not have to provide as many samples. In our use case Wage and Tax documents have a variety of formats and are not static.

- _6. Click the **ADD** button.

You should now see your new document type (class) in the list of classes on the left.

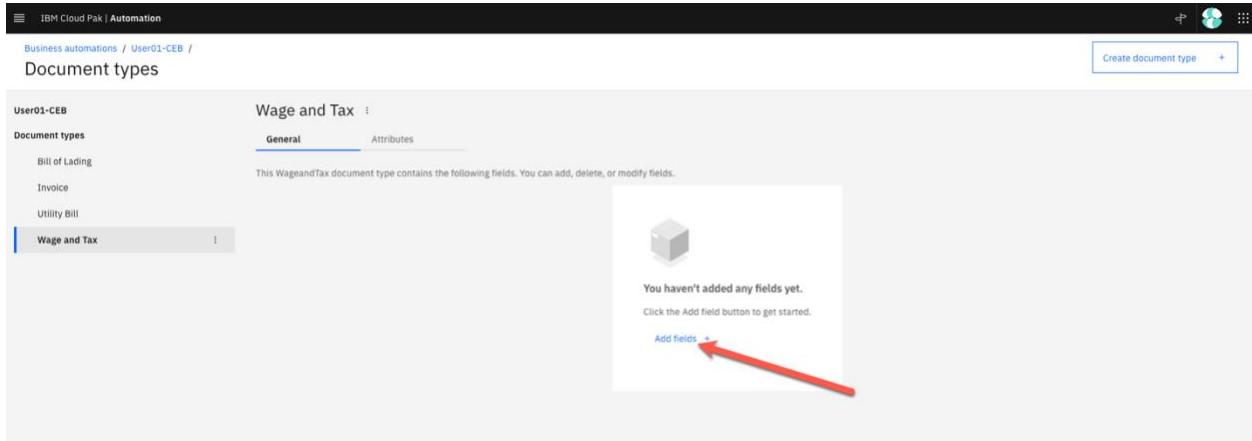


_7. Select your **Wage and Tax doc type**. On the right, you should see an empty table of fields.

5.2 Create Field

We can now add some fields to the class. From examination of the forms, we can see there are different fields names, or they are not consistent across the forms. We'll need to add these different "aliases" during this process.

_1. Click ADD FIELDS



_2. Enter the following values under the GENERAL Settings header

- Display Name: **Federal Income Tax Withheld**
- Field Type:
 - **Sys:Decimal**
- Is this field required: **Yes**
- In Aliases enter other possible names. Case and punctuation are very import when creating aliases. Enter the alias listed below. These are representations of what it looks like on the different forms. **Press the “+” after entering each one or press Enter key:**
 - **2 Federal income tax withheld**
 - **2. Federal income tax**



Note: the number two has a period after it

You should now see the following:

_3. Click the **NEXT** button.



Field patterns are regular expressions that can be associated with a field to help identify and extract fields. A regular expression is a sequence of characters that define a search pattern. The use of regular expression patterns and extractors is optional. Regular expression patterns can provide extra information to potentially improve the accuracy in extracting the correct fields. Python syntax is used for defining the regular expressions. You will not be adding any field patterns in this lab.

_4. Click **NEXT** again on the Field patterns screen. You should now be on the **VALUE SETTINGS** page. This is where you can set up validators, formatters, and converters.



Value Settings for a specific field; if the potential values follow a rule that can be expressed in a regular expression, you can specify an extractor. This pattern can match all the variations of your values. For example, the expected value for a Start Date field might be in a date format. You can create a regular expression pattern for `US Date` and then associate the extractor of `US Date` to your field.

Also, sometimes you want to extract a value that does not have a corresponding key in the document, but you know the pattern of the value. You can define the extractor and denote that the value might be anywhere in the document without attaching to the field name. This designation allows for the presence of a field name to be optional. For example, you want to extract the employee ID number, which can be described with a regular expression pattern. However, some documents show the employee number with a field name Employee ID, while other documents show the employee number without a corresponding field. You can specify the Extractor and be able to extract the employee ID number in both types of documents.

_5. The decimal data type can contain only integers to the left and right of a decimal point. But some of our data may contain commas between the integers and we only need two integers after the decimal point. Let's add a converter that will remove all extra punctuation and limit the number of integers after the decimal point to two. Click on the **Edit** button.

The screenshot shows the 'Value settings' tab in the 'Value format' section. It includes sections for 'Extractors', 'Formatters', and 'Converters'. A red box highlights the 'Edit' button next to the 'Converters' section.

_6. Click on Converters tab then click on blue button Add converter

The screenshot shows the 'Converters' tab selected in the 'Add field enrichments' section. A red box highlights the 'Converters(0)' tab.

_7. You will be presented with the Add converter screen. Click on Select existing. This populates the converter name, description, Decimal point, and Max digits after decimal point for you. If you wanted to change the decimal point from a period to a comma you could do it here as they do in other countries outside the United States. Click the blue Add button.

Add field enrichments to help the system extract the right data and reformat extracted values that might differ between documents. [Learn more](#)

Extractors(0) Formatters(0) Converters(0)

Add converter

How do you want to create a new converter?

Create new Select existing

Converter

Decimal Converter

Converter name
Decimal Converter

Description (optional)
Decimal Converter

Decimal point
.

Max digits after decimal point
2

[Cancel](#) [Add](#)

- _8. You will then be presented with Converter details information screen. On this screen you can also test your converters to make sure they are behaving like you intended. **Click on Done** at the top right.

Extractors(0) Formatters(0) Converters(1)

Converter details

Converter name
Decimal Converter

Description
Decimal Converter

Type
Decimal Converter

Decimal point
.

Max digits after decimal point
2

Inherited from
sys.Decimal

Test all converters

Sample value
Enter the sample value to test the converters with

[Test](#)

Converted result

No results yet
Converted result will be displayed here after clicking the Test button.

_ 9. Click **Create** your screen should look like this with your first field created.

Once it is created you will go back to the Document type page.

5.3 Create the Employee Name Address field.

_1. Click Add fields.

Give it the following parameters:

- Display name: **Employee Name and Address**
- Field Type = **sys:String**
- Required = **yes**
- Enter the following other possible names (aliases):
 - ***Employee name and address***
 - ***e Employee's first name and initial Last name Suff***
 - ***e Employee's name, address, and ZIP code***
 - ***e/f Employee's name, address, and ZIP code***
 - ***e. Employee Name & Address***
 - ***e Employee's first name and initial***

By default, the system will use the field name as an alias. So, you do not have to add it.

For example, below, Employee Name and Address (field name), would be automatically used as an alias even if you do not add it to the list

- _2. Click Next** no field patterns will be created.
- _3. Click Next** no value settings will be created.
- _4. Click Create** to finish creating the Employee Name and Address.

5.4 Create Employee Social Security Number Field

- _1. Click on ADD FIELDS**



Enter the following values in the GENERAL page.

- Display Name: **Employee Social Security Number**
- Field Type: **sys:Social Security Number**
- Is value required: **Yes**
- Other possible names (aliases). Remember, press RETURN or hit the '+' button on your keyboard between each entry:
 - **a Employee's social security number**
 - **a Employee's social security no.**
 - **a Employee's SSA number**
 - **a. Employee Social Security Number**
 - **Employee social security number**

Your screen should now look like the image below:

- _2. Click NEXT**
- _3. Click NEXT** again on the Field Patterns screen.
- _4. Click Create** on the Value settings.
- _5. Create the following additional Fields.**

The following table contains the values to use when adding the additional fields.

Follow the steps from the previous section to add the following fields. Don't forget to add your converter for datatypes of Sys:Decimal.

Display Name	Description	Type	Mandatory	Aliases
Employer Identification Number		sys:String	N	<ul style="list-style-type: none"> • b Employer identification number (EIN) • b Employer's FED ID number • b. Employer ID number • Employer identification number
Employers Name and Address		sys:String	N	<ul style="list-style-type: none"> • c Employer's name, address, and ZIP code • c Employer's Name & Address • Employers name and address
Social Security Wages		sys:Decimal	N	<ul style="list-style-type: none"> • Social security wages • 3 Social security wages
Wages Tips Other Compensation		Sys:Decimal	N	<ul style="list-style-type: none"> • 1 Wages, tips, other compensation • Wages, tips, other comp. • 1 Wages, tips, other comp. • 1. Wages tips, other comp • Wages tips other compensation

Reference for various field types:



Note: The basic default field types included in ADP are found here in the documentation

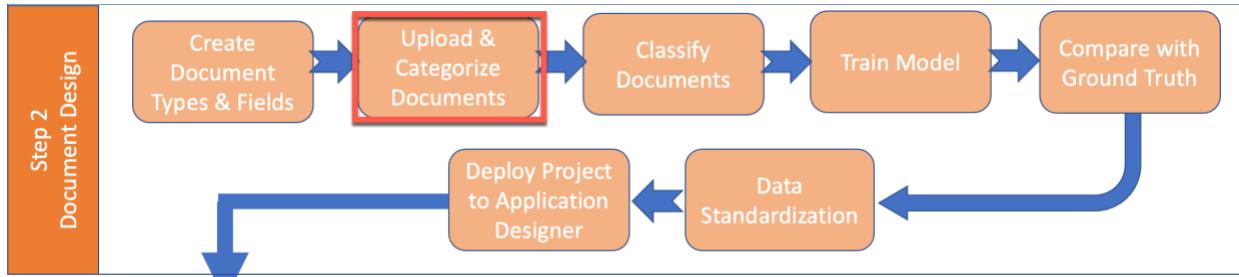
<https://www.ibm.com/docs/en/cloud-paks/cp-biz-automation/22.0.1?topic=enrichments-field-types-document-processing>

- _6. Click on the <name of your project> in the breadcrumb link in the top left of your screen. In the following example the name of the project is <Clandis Baker Project>. This will take you back to the Enrich tab, Then Click the Build tab.

The screenshot shows the IBM Cloud Pak Administration interface. The top navigation bar includes 'IBM Cloud Pak | Administration', a search bar, and a user icon. The breadcrumb path 'Business automations / Clandis Baker Project / Document types' is visible. On the right, there's a 'Create document type' button. The main content area is titled 'Wage and Tax' under 'Document types'. It has tabs for 'General' (selected) and 'Attributes'. A note states: 'This Wage and Tax document type contains the following fields. You can add, delete, or modify fields.' Below is a table of fields:

Name	Type	Required	Sensitive	⋮
Employee Name and Address	String	true	false	⋮
Employee Social Security Number	SocialSecurityNumber	true	false	⋮
Employer Identification Number	String	false	false	⋮
Employers Name and Address	String	false	false	⋮
Federal Income Tax Withheld	Decimal	false	false	⋮
Social Security Wages	Decimal	false	false	⋮
Wages Tips Other Compensation	Decimal	false	false	⋮

6 Document Types and Samples Overview



At this point in the process, we have created a new document type and configured the field names we want to extract off the document. For the system to know what to extract from your documents, it needs to be able to classify the documents. In this lab, we will teach the system to recognize the various document types on your system.

In the first part of the classification lab, you will explore the system's ability to automatically group similar documents together. This can be used to discover document types in a file share for example. You can also upload documents and have the system tell you what it finds. You would then use this information to create document types so you can classify the documents and data extract fields.

The project template comes pre-loaded with three document types: Bill of Lading, Invoice, and Utility Bill. In the last step we added a new document type Wages and Tax. In the BUILD tab of your project, you should now be seeing 4 document types. The three pre-loaded documents already have documents in them. You will be adding documents to the Wage and Tax document type. Your actual screen may vary from the following screen shot.

You will be asked to review the document categories the system finds and create the appropriate document types as needed.

6.1 Categorize documents.

For categorizing, we will have the system help us group similar documents together. To get started,

1. Click anywhere in the Document types and samples box.

The screenshot shows the 'Build' tab selected in the navigation bar. The main content area displays several sections: 'Document types and samples' (highlighted with a red box), 'Classification model', 'Extraction model', and 'Data standardization'. The 'Document types and samples' section includes a sub-instruction: 'Upload sample documents to define the types of documents you want the system to process.' To the right of this section, there are status indicators: 'Ready' (with a green dot), '4 types', '22 samples on average', and '100% accuracy'. Below this, the 'Classification model' section shows 'Ready' (green dot), '3 types trained', and '100% accuracy'. The 'Extraction model' section shows 'Ready' (green dot), '3 types trained', and '97% accuracy'. The 'Data standardization' section shows 'Not ready' (yellow dot).

The CATEGORIZE feature analyzes each document and tries to find similarities between them. Based on these similarities, the system will divide the samples into categories for you to review. You can add documents or entire categories into either an existing document class or create new classes as needed.

Let's see what that looks like.

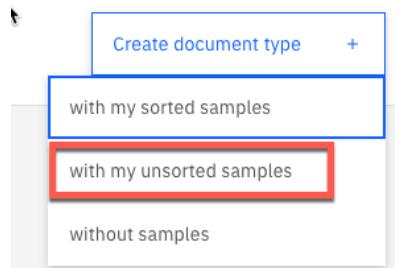
2. Click on **CREATE DOCUMENT TYPE** in the top right of the screen.

The screenshot shows a dropdown menu triggered by a button labeled 'Create document type +'. The menu contains three options: 'with my sorted samples' (highlighted with a blue box), 'with my unsorted samples', and 'without samples'. There is also a small 'Upload' button at the bottom of the menu.

If you have the same document types already separated into folders, you can choose the first option, *with my sorted samples*. The system would simply ingest the documents from each folder into a different group.

For this exercise, we will select the second option, *with my unsorted samples* and let the system sort the documents for us. Use this option when you don't know how many different document types there are.

_3. Select the second option titled **with my unsorted samples.**



You should have already downloaded the files from [Section 2](#) to your laptop. You can select upload and grab all the files from where they were downloaded to on your laptop. Make sure you have already unzipped them.

_4. Click Upload to get document samples.

From the downloaded sample documents open the folder name [Group 1 – Design Docs for Tax Lab.](#)

Note: this will take several minutes, good time for some coffee or a stretch. Make sure to check ALL documents have been uploaded there are two pages or 12 items to verify.

At the bottom of the window, you can select the number of items to display in the window or click on the arrows to move to the next page

Upload sample documents that represent the different types of documents you want the system to classify. Include at least 6 samples of each type of document.

Search sample documents Upload

Document name	
<input type="checkbox"/> Mortgage Agreement1.pdf	✓
<input type="checkbox"/> Mortgage Agreement2.pdf	✓
<input type="checkbox"/> Mortgage Agreement3.pdf	✓
<input type="checkbox"/> Mortgage Agreement4.pdf	✓
<input type="checkbox"/> Mortgage Agreement5.pdf	✓
<input type="checkbox"/> TR_FW2_1001_0000_P5.pdf	✓
<input type="checkbox"/> TR_FW2_2000_0000_P5.pdf	✓
<input type="checkbox"/> TR_FW2_3000_0000_P5.pdf	✓
<input type="checkbox"/> TR_FW2_3001_0000_P5.pdf	✓
<input type="checkbox"/> TR_FW2_4000_0000_P5.pdf	✓
<input type="checkbox"/> UBILLCable_081_1_1.1.pdf	□
<input type="checkbox"/> UBILLCable_082_1_1.1.pdf	□

Items per page: 20 ▾ 1 - 12 of 12 Items 1 of 1 page ▶ ▷

5. Click on the blue CATEGORIZE button on the top right corner.



Note: The results may vary based on the documents uploaded, what the system already has learned, the version of ADP and more. Please look at this lab exercise from a high level. The categories you will be presented are the system's best guess on how they should be separated.

You will need to:

- Review the categories to see if the documents were separated correctly.
- Move documents into either a NEW document type or into an EXISTING document type.
- There should be 3 types in the samples you were provided.
 - Wage and Tax
 - Utility bills
 - Mortgage Agreements
- You will need to assign either an entire category (i.e., all sample documents) or individual documents in each category to the Wage and Tax and Utility bills document types which already exist on your system.
- You will need to create a new document type for Mortgage Agreements.

After a few seconds, the system will mark the documents with a status of ready as seen in the above image.

6. Click on each of the categories to see what was grouped together as shown below.

You can Click on any document to see a preview of it. This will help ensure the documents are correctly grouped.



NOTE: The names of the files are not used in any way in this process. The files were merely named this way to make it easier for you to quickly ascertain whether the documents were grouped correctly.

This screenshot shows the 'Create document types' interface in IBM Cloud Pak Automation. The left sidebar lists 'Categories (3)': Category 1 (selected), Category 2, and Category 3. Under 'Document types (4)': Bill of Lading, Invoice, Utility Bill, and Wage and Tax. The main area shows 'Category 1 sample documents (2)'. A search bar 'Search sample documents' and an 'Upload' button are at the top. Below is a list of documents: UBILLCable_081_1_1.1.pdf and UBILLCable_082_1_1.1.pdf, both with checkmarks.

This screenshot shows the 'Create document types' interface in IBM Cloud Pak Automation. The left sidebar lists 'Categories (3)': Category 1 (selected), Category 2, and Category 3. Under 'Document types (4)': Bill of Lading, Invoice, Utility Bill, and Wage and Tax. The main area shows 'Category 2 sample documents (5)'. A search bar 'Search sample documents' and an 'Upload' button are at the top. Below is a list of documents: Mortgage Agreement1.pdf, Mortgage Agreement2.pdf, Mortgage Agreement3.pdf, Mortgage Agreement4.pdf, and Mortgage Agreement5.pdf, all with checkmarks.

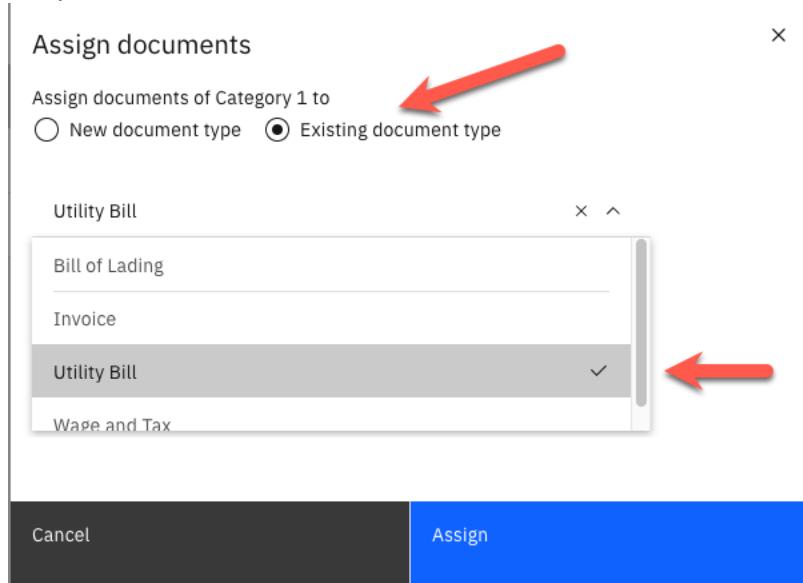


At the time of writing this documentation ADP was able to categorize the sample set into each category. This is not always the case, sometimes document types will be combined into one category, so it's very important to look at each category and verify documents.

- _7. If all documents within a category are correct as illustrated in the following screen shot, **Click** on the **3 dots** at the end of the category name.

- _8. Select Assign to document type**

_9. Select Existing Document type then the appropriate document type from the drop-down list.

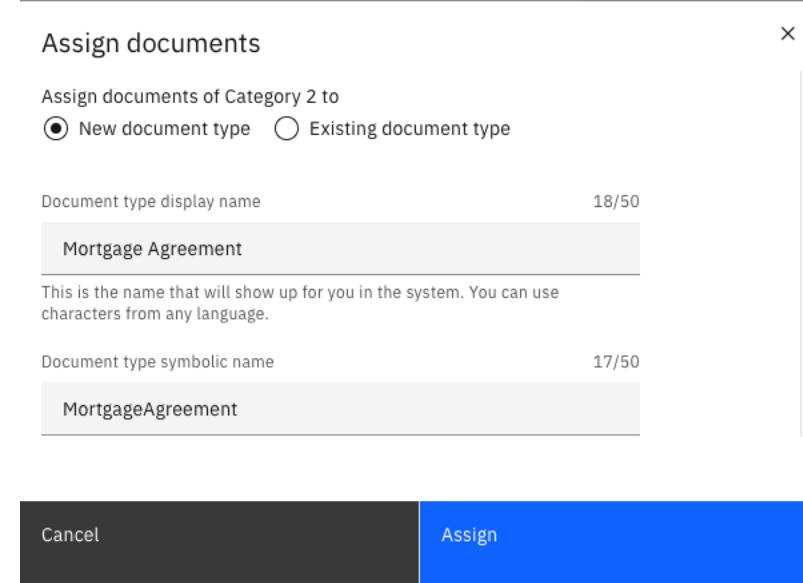


_10. Click Assign to close the dialog box.

_11. Select the next Category 2 and Click on the 3 dots and Select Assign to document type to a document class.

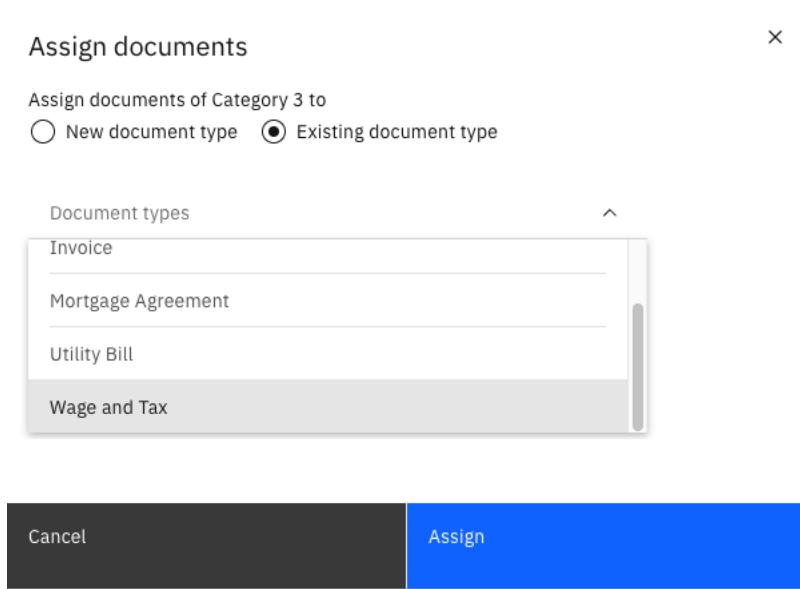
_12. This time Select a New Document Type. Since we have not defined a mortgage agreement document type yet.

_13. Enter Mortgage Agreement in the field



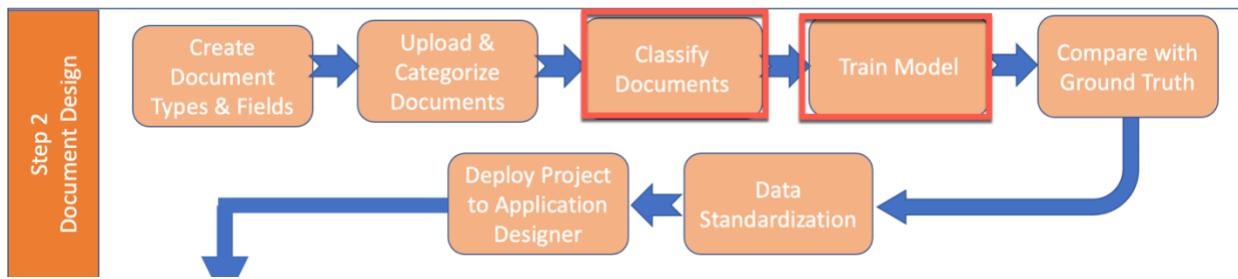
_14. Click Assign to have the system automatically rename and move the category into the Document Types section.

- _15. Now for Category 3, **Click on 3 dots** and **Select Assign document type**.
- _16. **Click Existing document type** and **Select Wage and Tax** from the drop down and then **Click on Assign**.



- _17. Once you confirmed all documents are correctly classified into the correct document type, **Click Finish**

7 Train classification



Now that we have documents uploaded in the system, we are ready to train the classification. Note that although you don't need a ton of document samples to train (minimum of 5), you are going to get better accuracy if the system has a deeper understanding of the documents, so more could be better.

In this lab, we curated some documents samples for you. In normal circumstances, you would need to do this yourself. Make sure the documents you upload to train classification are good documents.

- Clean documents
- High resolution
- Representative of the document type(s)
- Accurately grouped and uploaded to Document Processing

This is NOT the time to try and trick the system. Uploading a document that doesn't recognize well would not help the system recognize the types of words, phrases, and concepts it needs to learn to classify documents correctly.

The most common error is introducing a sample document into the incorrect document type, usually by uploading them to the wrong document type. If that happens, you are introducing conflict into the classification. For example, an invoice added to Tax Forms may confuse the system and result in it thinking invoices are tax forms and vice versa. Once that happens, you need to clean your documents and retrain the system.

- _1. Click on <your project name>** in the cookie trail to return to the start page.
- _2. Click anywhere in the **CLASSIFICATION MODEL** line**

Document types and samples

Upload sample documents to define the types of documents you want the system to process.

	Status	Types Trained	Accuracy
Classification model	Ready	3	100% accuracy
Extraction model	Retrain	3	97% accuracy
Data standardization	Not ready		
Document retention	Ready	5	types reviewed

Once we open the classification model, we will be presented with details on how to perform the retraining. There are four basic steps – Confirm inputs, Review Samples, Review Training Results, and Test Trained model.

On the Confirm inputs screen here we can confirm all the documents that will be used in this training exercise. We can also use the opportunity to remove documents that are no longer relevant or upload additional documents.

- _3. **Click Next** this will move from the Confirm inputs to the **Review Samples** step. Notice three documents have green icons next to them. These green icons show these documents have test samples already assigned. The new document types (Mortgage Agreement and Wage and Tax) do not have any test samples assigned yet therefore there's no green icons since we haven't assigned test sets yet.

Classification model

Accuracy: 84.8%

Document types:

- Bill of Lading
- Invoice
- Mortgage Agreement
- Utility Bill
- Wage and Tax

Mortgage Agreement sample documents (5)

Training/test ratio in %: 100/0

Test set (0)

0% of total samples

There are no documents in the test set.

_4. For the Mortgage Agreement move two documents to the Test set by **checking** and **click** on the **arrow** in between columns.

Classification model

Accuracy: 84.8%

Document types:

- Bill of Lading
- Invoice
- Mortgage Agreement
- Utility Bill
- Wage and Tax

Mortgage Agreement sample documents (5)

Training/test ratio in %: 60/40

Test set (2)

40% of total samples

_5. Select **Wage and Tax** on the Document types. This time let the ADP system Auto generate the 70/30 split to the test set. **Click Auto generate 70/30 split**

Classification model
Last trained: 5 months ago
Accuracy: 100%

Confirm inputs Review samples Review training results Test trained model Optional

Changes were made since you last trained your model. Retrain the model to get updated training results and accuracy.

Document types

- Bill of Lading: 28 samples
- Invoice: 24 samples
- Mortgage Agreement: 5 samples
- Utility Bill: 35 samples
- Wage and Tax**: 5 samples

This document type will not be trained because you have no documents in the test set. Please make sure you have at least 1 document in each set.

Review your training and test sets. A good practice is to assign 70% of your samples to the training set and 30% to the test set. The test set is used to generate the model training results. [Learn more](#)

Wage and Tax sample documents (5) Training/test ratio in %: 100/0

Training set (5) 100% of total samples

Search training set sample documents

- TR_FW2_1001_0000_P5.pdf
- TR_FW2_2000_0000_P5.pdf
- TR_FW2_3000_0000_P5.pdf
- TR_FW2_3001_0000_P5.pdf
- TR_FW2_4000_0000_P5.pdf

Test set (0) 0% of total samples

Search test set sample documents

There are no documents in the test set. Include at least 1 document in the test set to view training results.

Auto generate 70/30 split



The suggested split is 70/30 – that is, 70% of the available sample documents should be used for training, and we will validate the training results with 30% of the sample documents. This split is only a suggestion, and we can adjust it, but 70/30 is a good starting point.

Classification model
Last trained: a day ago
Accuracy: 84.8%

Confirm inputs Review samples Review training results Test trained model Optional

Changes were made since you last trained your model. Retrain the model to get updated training results and accuracy.

Document types

- Bill of Lading: 23 samples
- Invoice: 31 samples
- Mortgage Agreement: 5 samples
- Utility Bill: 27 samples
- Wage and Tax**: 5 samples

Review your training and test sets. A good practice is to assign 70% of your samples to the training set and 30% to the test set. The test set is used to generate the model training results. [Learn more](#)

Wage and Tax sample documents (5) Training/test ratio in %: 60/40

Training set (3) 60% of total samples

Search training set sample documents

- TR_FW2_3000_0000_P5.pdf
- TR_FW2_3001_0000_P5.pdf
- TR_FW2_4000_0000_P5.pdf

Test set (2) 40% of total samples

Search test set sample documents

- TR_FW2_1001_0000_P5.pdf
- TR_FW2_2000_0000_P5.pdf

Auto generate 70/30 split

6. Click on TRAIN to launch the training. This may take a several minutes. You will see a progress bar has training progresses.

Once complete, you will be able to see the training results.



What's happening: All the samples are run through multiple machine learning algorithms. These machine learning algorithms learn from the ground truth, the association between the sample documents (the OCR text) and the document types. The yielding models are then evaluated with the documents in test set. The model-predicted document types on these documents are compared with the human-provided answers to compute the accuracy. The top three accurate models are presented to the user, with the most accurate one being selected by default.

You should see something like the following:

Document	Classified as	Classification result	Confidence
BOL_007.pdf	Bill of Lading	Correct	High
BOL_009.pdf	Bill of Lading	Correct	Medium
BOL_019.pdf	Bill of Lading	Correct	High
BOL_027.pdf	Bill of Lading	Correct	High
BOL_031.pdf	Bill of Lading	Correct	High
BOL_075.pdf	Bill of Lading	Correct	High

7. Click on each of the document types. Notice the confidence levels. The both the Mortgage Agreement and Wage and Tax have a confidence of low. Low Confidence means we probably need to add more documents to our document class to get better confidence values.



You can easily see where the system may be struggling with Wage and Tax and Mortgage Agreement. You should look for document types that don't match the actual file or have a low confidence. Remember the more documents you give to train, the better the results.

_8. **Click on Next.** This is the Test trained model. Here you can try and test other documents to see if they classified correctly. This step is optional but would be useful to try out the AI model to determine whether additional samples are necessary.

_9. **Click Done**

7.1 How do I improve my results?

7.1.1 Option 1 – Add more samples.

To improve results, you would normally want to add more samples of the document ensuring they are clean and representative document to improve the system's understanding of the document.

_1. **Click anywhere on Document Types and Samples.**

_2. **Click on Wage and Tax type.**

_3. **Click on Upload**

- _4. From the zip files you downloaded and unzipped earlier upload all the files from the directory *Group 2 - Classification Results Increase Set*.
- _5. Go back to the **Build** tab then let's retrain the Classification Module again.
- _6. **Click anywhere on Classification model.**
- _7. **Click on Wage and Tax.**

The screenshot shows the 'Document types and samples' section of the IBM Cloud Pak Administration interface. On the left, there is a sidebar with a list of document types: Bill of Lading (28 samples), Invoice (26 samples), Mortgage Agreement (5 samples), Utility Bill (35 samples), and Wage and Tax (5 samples). The 'Wage and Tax' item is selected and highlighted with a blue border. To the right, there is a main panel titled 'Wage and Tax sample documents (5)'. It contains a search bar labeled 'Search sample documents' and a list of five PDF files, each with a checkbox and a checkmark icon. At the bottom right of this panel, there is a blue 'Upload' button with a white arrow icon, which is also highlighted with a red box. The top navigation bar includes links for 'IBM Cloud Pak | Administration', 'Business automations / Claudi Baker Project /', and 'Create document type +'. There are also icons for user profile and settings.

_8. Click Next button. Also click on the Auto generate 70/30 split.

The screenshot shows the 'Classification model' page in the IBM Cloud Pak Administration interface. The 'Wage and Tax' document type is selected. The 'Training set (7)' contains 70% of total samples, and the 'Test set (3)' contains 30% of total samples. A red box highlights the 'Auto generate 70/30 split' button in the top right corner of the main content area.

_9. Click Train button.

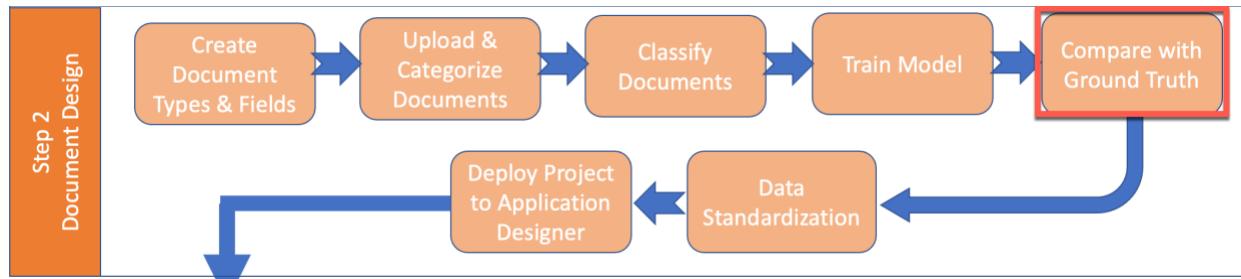
_10. Now look at the confidence score for **Wage and Tax**.

_11. Click Next and then Click Done

7.1.2 Option 2 – Review all uploaded samples.

- remove those that are not a clear representation.
- remove those that are poor quality documents.
- carefully confirm that none of the samples contain multiple document types in the file. This is a common occurrence. A document is listed as a Purchase Order, but in the back pages, also contains other document types in that same file. This confuses the system.

8 Data extraction



At this point, we have defined a document type, told the system which fields we want off the document and trained the system on how to recognize (classify) the document. In the Data Extraction portion of the lab, we will upload new Wage and Tax documents to Document Processing and see how our earlier configuration of the document type and related fields are working. This is comparing a new document extracted elements with the ground truth. Once we open Extraction model, we will be presented with details on how to perform the retraining. There are five basic steps – Review samples, Add fields, Teach the model, Review the trained model, and Test the model.

- _1. From the guided configuration screen, **Click** anywhere in the **Extraction model** box.



Note: the status will reset to Retrain if it detects something may have changed. This is just a reminder that if you indeed changed something, you may benefit from retraining the model.

Extraction model

Train the model to extract the data from your documents.

! Retrain

3 types trained
90% accuracy

Open →

- _2. Next **Click** on the **Wage and Tax** document type under the Document Types section.

Like in the classification step, ADP needs to have the documents divided into a training and test sets. In general, *deep learning*-based AI requires a larger number of sample documents to achieve a reasonable result. But since our environment does not have GPU, deep learning is not turned on.

You should have something that looks like what you see in the following screen shot.

_3. Again, lets train with a 70/30 spilt. **Click Auto generate 70/30 split.**

_4. **Click** on the **NEXT** button at the top.



You will now be on the Add fields bread crumb. If there were more fields to add we could do it here. But since we have already added all the fields needed, proceed to the next step.

_5. **Click** the **Next** button. You are now at the “Teach model” bread crumb.

Teach the model is where you will spend most of your time. We can see that our documents are “not ready”, so we’ll need to teach the model with new documents.

_ 6. Click on Teach Samples.

The screenshot shows the 'Teach model' step in the IBM Cloud Pak Administration interface. On the left, there's a sidebar with 'Document types' including Bill of Lading, Invoice, Mortgage Agreement, Utility Bill, and Wage and Tax. The 'Wage and Tax' section is selected and highlighted with a blue border. In the main area, a yellow warning box says 'Please make sure you have at least 1 reviewed document to train the model.' Below it, a list of 'Wage and Tax sample documents (7)' is shown, each with a status of 'Not ready' and 0/7 fields reviewed. At the bottom right of the main area, there are 'Upload' and 'Reanalyze' buttons, and a red box highlights the 'Teach samples' button.



Note: Your individual results may vary based on the exact documents you upload, how you configure your fields etc. Therefore, general guidance is given here versus exact step by step instructions.

_ 7. We will now review the fields that were extracted, correct any that may be wrong and add others.

You should now see the field data extracted by the system. Nothing has been trained yet. All it is doing is using the field name and aliases we entered when we created the document class to locate data. Now, you need to correct and improve the model.

The screenshot shows the IBM Cloud Pak Administration interface. On the left, a PDF of a W-2 form for the year 2020 is displayed. The form includes fields for employer information, employee details, and various tax withholdings. On the right, a sidebar titled "Match data underlined in blue to the selected field or draw your own boxes around data in the document." contains a table titled "Field Name Value Captured". The table lists several extracted fields with their corresponding values. A "Draw" button is available for each entry to refine the extraction. At the bottom of the sidebar, there are buttons for "Previous sample" and "Next sample".

Field Name	Value Captured
Federal Income Tax Withheld	abc Text Required
Field label (optional)	Draw <input type="text"/> Captured field label
Field value	Draw <input type="text"/> Captured field value
Pending aliases	View all aliases (3)
None ⓘ	
Save selection	
Employee Name and Address	abc Text Required
Employee Social Security Number	abc Text Required
Employer Identification Number	abc Text
Employers Name and Address	abc Text
Social Security Wages	abc Text
<input type="checkbox"/> Mark this document as ready for training. ⓘ	
Previous sample	Next sample



Note: You may see different results than shown on the image above. Depending on how the algorithms interpreted the results you could see either type of extraction.

This screenshot shows the same W-2 form and extraction interface as the previous one, but with different results. The sidebar now displays a "Recommended matches" section, which lists "Federal income tax withheld" with a value of "1800.00". This indicates that the system has identified a different set of data as being more likely to be correct. The rest of the sidebar and table structure are identical to the first screenshot.

Field Name	Value Captured
Federal Income Tax Wit...	abc Text
Field label	Field value
Federal income tax withheld	1800.00
2 Federal income tax withheld	1800.00
<input type="button"/> Edit selection <input type="button"/> Dismiss <input type="button"/> Seeing duplicates? ⓘ	
Pending aliases	View all aliases (3)
Detected alias already exists ⓘ	
Save selection	
Employee Name and A...	abc Text
Employee Social Securi...	abc Text

Let's spend some time showing how to go about correcting these issues to help the system learn how to extract the values accurately.

8.1 Correcting extracted values

Let's start with the Federal Income Tax withheld field (i.e., The first one in the 'Fields to extract' list). Again, you may see different results based on your forms and how the different algorithms behaved on that particular document during extraction.

_1. ADP may have already preselected the first field like in the first screen shot below.

But ADP can also show the characters it recognized on the page with blue lines (second screen shot below) If your result is like the first screen shot then **Click** blue button **Save section**. Otherwise, if you got blue lines **Click** on the **number** below the heading "**Federal Income tax withheld**" in the image.

Field Name	Value Captured
Federal Income Tax Wit...	1800.00

Recommended matches

- Federal income tax withheld 1800.00
- 2 Federal income tax 1800.00 withheld
- Federal income tax 1800.00 withheld
- Local income tax 500.00

Action Buttons: Edit selection, Dismiss, Seeing duplicates?

W-2 Wage and Tax Statement
Form **W-2** Wage and Tax Statement
Copy 1—For State, City, or Local Tax Department

2020

Department of the Treasury—Internal Revenue Service

IBM Cloud Pak | Administration

← Back TR_FW2_1001_0000_PS.pdf | Not ready

Show detected fields Keyboard shortcuts on

577-22-3048 OMB No. 1545-0008

22222	a Employee's social security number	577-22-3048	OMB No. 1545-0008
b Employer identification number (EIN)	14-023285		
c Employer's name, address, and ZIP code	Long Lengthy Name The Corporation 56334 Full Sized Avenue Unit 1234 Minneapolis, Minnesota 55411-1234		
d Control number	123456 A78		
e Employee's first name and initial last name	Michael Robert David Smithson III 56334 Full Sized Avenue Unit 1234 Minneapolis, Minnesota 55411-1234		
f Employee's address and ZIP code			
15 State	Employer's state ID number	16 State wages, tps, etc.	17 State income tax
MN	123456789	123456789.99	123456789.99
18 Local wages, tps, etc.	19 Local income tax	20 Locality name	
123456789.99	123456789.99	ABCDEF	

W-2 Wage and Tax Statement 2020 Department of the Treasury—Internal Revenue Service

Copy 1—For State, City, or Local Tax Department

Sort by: Date created

Field Name Value Captured

- Federal Income Tax Withheld Required
- Field label (optional) Draw Captured field label
- Field value Draw Captured field value
- Pending aliases | View all aliases (3) None
- Save selection

Employee Name and Address Required

Employee Social Security Number Required

Employer Identification Number

Employers Name and Address

Mark this document as ready for training.

Previous sample Next sample

- _2. Again, depending on your specific results. If ADP was able to find the field and will ask if you want to save match of value captured along with the field label. **Select Save Selection**. Otherwise, If your results were the recognized characters with blue lines then in the pop up window that comes up **select Save match**

IBM Cloud Pak | Administration

← Back TR_FW2_1001_0000_PS.pdf | Not ready

Show detected fields Keyboard shortcuts on

577-22-3048 OMB No. 1545-0008

22222	a Employee's social security number	577-22-3048	OMB No. 1545-0008
b Employer identification number (EIN)	14-023285		
c Employer's name, address, and ZIP code	Long Lengthy Name The Corporation 56334 Full Sized Avenue Unit 1234 Minneapolis, Minnesota 55411-1234		
d Control number	123456 A78		
e Employee's first name and initial last name	Michael Robert David Smithson III 56334 Full Sized Avenue Unit 1234 Minneapolis, Minnesota 55411-1234		
f Employee's address and ZIP code			
15 State	Employer's state ID number	16 State wages, tps, etc.	17 State income tax
MN	123456789	123456789.99	123456789.99
18 Local wages, tps, etc.	19 Local income tax	20 Locality name	
123456789.99	123456789.99	ABCDEF	

W-2 Wage and Tax Statement 2020 Department of the Treasury—Internal Revenue Service

Copy 1—For State, City, or Local Tax Department

Sort by: Date created

Field Name Value Captured

- Federal Income Tax Withheld Required
- Field label (optional) Draw 2 Federal income tax withheld
- Field value Draw 123456789.99
- Pending aliases | View all aliases (3) None
- Save selection

Employee Name and Address Required

Employee Social Security Number Required

Employer Identification Number

Employers Name and Address

Mark this document as ready for training.

Previous sample Next sample

Notice a green check mark signifies this field is complete.

The screenshot shows the IBM Cloud Pak Administration interface. On the left is a PDF of a W-2 form for the year 2020. The form includes fields for employer information, employee name and address, and various tax components like wages, tips, and taxes withheld. On the right, a sidebar displays the captured data from the form. A specific field, "Federal Income Tax Withheld" (value 123456789.99), has a red box drawn around it. Below this, there are three ellipses (three dots) next to a green checkmark, which are used for clearing or updating data.

Field Name	Value Captured
Federal Income Tax Withheld	123456789.99
Field label (optional)	2 Federal income tax withheld
Field value	123456789.99

The 3 ellipses next the green check mark allow you to clear the data or update ADP to there is no field with this data in the current view.

- _3. Move to Employee Name and Address field by clicking in the grey area on that field name. In our two possible outcomes depending on the algorithms. ADP did pick up the address but missed the name. Or the algorithm may have picked up the address and not the name. Or it may have gotten the correct field.

If the field is not correct **Click on the Dismiss button**.

Now under the Field label **select Draw** button and using your mouse grab or lasso around “Employee’s first name and initial”.

This screenshot shows the same W-2 form and data capture interface as the previous one, but with a different focus. The "Employee's first name and initial" field (e) is now highlighted with a blue selection box. To the right, a "Recommended matches" panel lists several suggestions, including "e Employee's first name and initial" with the value "4326 Aldrich Rd Minneapolis, MN 55412". This indicates that while the name was initially missed, the system was able to correctly identify it during the document processing.

Field Name	Value Captured
e Employee's first name and initial	4326 Aldrich Rd Minneapolis, MN 55412

If you got the blue lines, you would notice that only the “Employee’s first name and initial” have blue marks. In this case the values for name and address were not located. Using Draw button and using your mouse grab or lasso around “Employee’s first name and initial”.

- _4. We are interested in getting the “Employee’s First Name” data and address for the field value. **Click** on the **Draw** button under Field value. Using your mouse select the appropriate values for Name and address (green box), then **Click Save selection**

The screenshot shows the IBM Cloud Pak Administration interface with the following details:

- W-2 Wage and Tax Statement Form:** The form is displayed with various fields filled in. A green box highlights the "Employee's first name and initial" field, which contains the value "Benjamin P. Charles".
- Extracted Data:** On the right, a table shows extracted data from the form. The "Employee Name and Address" row has a green box around it, containing the value "Benjamin P. Charles, 4326 Aldrich Rd, Minneapolis, MN 55412".
- UI Elements:** Several "Draw" buttons are highlighted with red boxes. A large red box surrounds the "Employee Name and Address" row in the extracted data table. A blue box highlights the "Save selection" button at the bottom right of the interface.

- _5. For the Employee Social Security field if it looks good, **Click** on **Save selection**. Or if the blue lines are present instead **select** the value displayed to populate the field and **Click Save match** then **Click on Save selection**.
- _6. Continue to process for the remaining fields, using either method as described above, clicking on the **Save selection** if ADP picked up the correct field label and field value or select the blue line values to populate both the field label and field value or finally if both fields are wrong use the **Dismiss** and use blue lines if Key Value Pair (KVP) is correct or drawing a box around needed label or value.

- _7. Once complete **check the box** next to “Mark this document as ready for training” at the bottom

The screenshot shows the IBM Cloud Pak Administration interface. On the left, there is a PDF preview of an W-2 Wage and Tax Statement form. The form includes fields for employee information (SSN: 577-22-3048, EIN: 14-023285), employer details (Test and Rest Inc., 563 Stoney Brook Rd, Minneapolis, MN 55411), and tax withholdings. On the right, the extracted data is shown in a structured format. A red arrow points to the checkbox labeled "Mark this document as ready for training".

22222	a Employee's social security number 577-22-3048	OMB No. 1545-0008		
b Employer identification number (EIN) 14-023285	1 Wages, tips, other compensation 18000.00			
c Employer's name, address, and ZIP code Test and Rest Inc. 563 Stoney Brook Rd Minneapolis, MN 55411	2 Federal income tax withheld 1800.00	3 Social security wages 17700.00		
d Control number 210220 A13	4 Social security tax withheld 1113.33	5 Medicare wages and tips 18000.00		
e Employee's first name and initial Benjamin P. Charles 4326 Aldrich Rd Minneapolis, MN 55412	6 Medicare tax withheld 261.00	7 Social security tips 400.00		
f Employee's address and ZIP code 15 State Employer's state ID number MN 795037	8 Allocated tips 400.00	9 Dependent care benefits 543.21		
	10	11 Nonqualified plans 300.00		
	12a	12b A 256.00		
	13	13b Employee Retirement Plan X X X 20000.00		
	14	14c Other Test form DD 532.00		
		14d AA 425.00		
16 State wages, tips, etc. 18000.00	17 State income tax 1260.00	18 Local wages, tips, etc. 17700.00	19 Local income tax 500.00	20 Locality name MPLS

W-2 Wage and Tax Statement
Form 1—For State, City, or Local Tax Department
2020
Department of the Treasury—Internal Revenue Service

- 8.**  Review ALL other fields carefully. **Do not leave any incorrect values.** You can adjust or delete values as needed by clicking on Edit selection. If you leave incorrect values, the system will assume they are correct and LEARN them as if they were good values.

9. Repeat steps for Next Sample

Over the course of next few samples you may find that ADP has extracted the wrong results, perhaps getting a value that is above when it should have been below. If this is the case and you pick you a blue underline data, but the results are wrong. Simply use the draw box for the Field Label and Field Value.



Note: When completing the remaining documents, you may run across ADP finding the fields but perhaps on the second image or third image on the page. Try to keep all Key Value Pairs (KVP) on the same image.

_ 10. Once complete review of all the sample documents **Click on the Back link**

The screenshot shows the IBM Cloud Pak Administration interface with three W-2 forms displayed side-by-side. Each form is a table with various tax-related fields. A red box highlights the 'Back' button at the top left of the first form. To the right of the forms is a search bar with the placeholder 'Search document, tax, VML...'. Below the search bar are several dropdown menus and input fields, some with validation messages like 'Required' or 'Field label (optional)'. At the bottom right, there are buttons for 'Save selection' and 'Mark this document as ready for training'.

Form W2 Wage and Tax Statement Copy 1 - For State, City or Local Tax Department OMB No. 1545-0008			
a. Employee Social Security Number: 328-47-1017	b. Employer ID number: 98-7459372	c. Control number: AA346 45	
d. Employee Name & Address David Gomez 563 Broadway New York, NY 10027			
1. Wages tips, other comp. 210000.00	2. Federal income tax 50400.00	3. Social security wages 132099.00	4. Social security tax withheld 83239.80
5. Medicare wages and tips 190000.00	6. Medicare tax withheld 2755.00	7. Social security tips Allocated tips	
9	10. Dependent care benefits	11. Nonqualified Plans	12a. Code (see instructions) C 1423.00
13. (Mark) d Statutory employee Retirement plan X	14. Other	12b. Code D 20000.00	12c. Code DD 532.00
15. State NY	State ID number 795037	16. State wages, tips, etc. 210000.00	17. Social security wages FF 425.00
17. State income tax 14700.00	18. Local wages, tips, etc.	19. Local income tax Red Beach	20. Locality Name

Form W2 Wage and Tax Statement Copy 2-To Be Filed with Employee's FEDERAL Tax Return OMB No. 1545-0008			
a. Employee Social Security Number: 328-47-1017	b. Employer ID number: 98-7459372	c. Control number: AA346 45	
d. Employee Name & Address David Gomez 563 Broadway New York, NY 10027			
1. Wages tips, other comp. 210000.00	2. Federal income tax 50400.00	3. Social security wages 132099.00	4. Social security tax withheld 83239.80
5. Medicare wages and tips 190000.00	6. Medicare tax withheld 2755.00	7. Social security tips Allocated tips	
9	10. Dependent care benefits	11. Nonqualified Plans	12a. Code (see instructions) C 1423.00
13. (Mark) e Statutes emulsione	14. Other	12b. Code D 20000.00	

Form W2 Wage and Tax Statement Copy C - For Employee's Records OMB No. 1545-0008			
a. Employee Social Security Number: 328-47-1017	b. Employer ID number: 98-7459372	c. Control number: AA346 45	
d. Employee Name & Address David Gomez 563 Broadway New York, NY 10027			
1. Wages tips, other comp. 210000.00	2. Federal income tax 50400.00	3. Social security wages 132099.00	4. Social security tax withheld 83239.80
5. Medicare wages and tips 190000.00	6. Medicare tax withheld 2755.00	7. Social security tips Allocated tips	
9	10. Dependent care benefits	11. Nonqualified Plans	12a. Code (see instructions) C 1423.00
13. (Mark) f Statutes emulsione	14. Other	12b. Code D 20000.00	

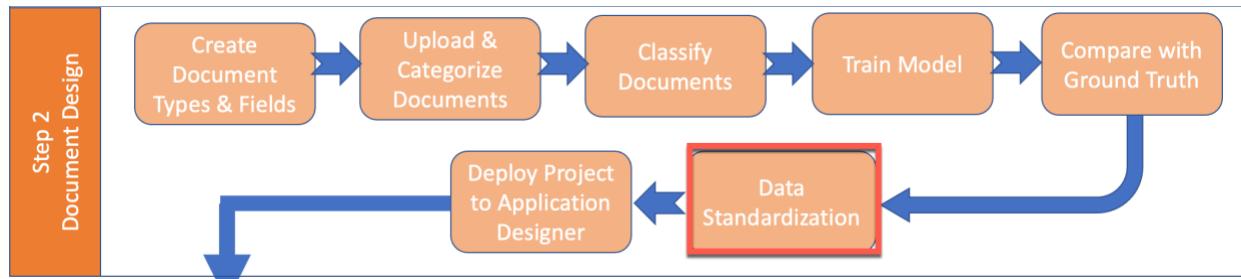
8.2 Train extraction model

We will be performing the quick training in this lab due not having a GPU in our TechZone architecture. A GPU is only needed a development environment and is not needed in either a production or runtime environment. The Deep Learning capabilities have been disabled on this training environment. You can find instructions in the Appendix for when you have access to a server with it enabled.

_ 1. Click Train button.

This will take several minutes. (Good time for a break)

9 Data standardization

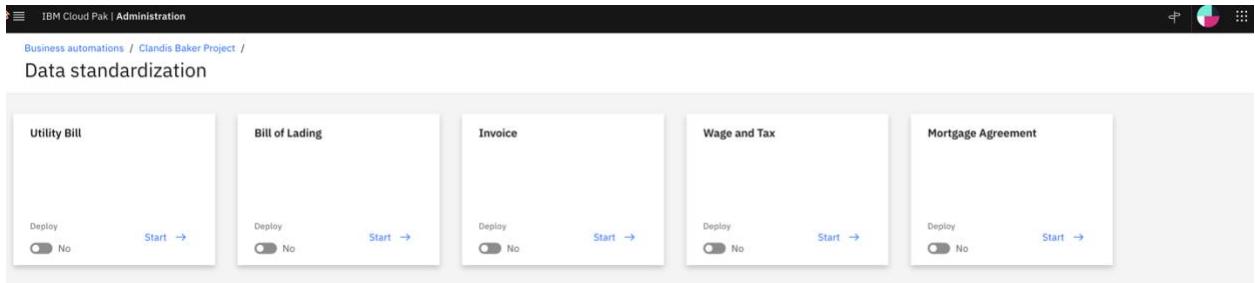


Next, we may need to standardize the data that will be presented in the user interface and how it will be stored in the FileNet repository for example. Data standardization is the process of defining attributes for a data field in a standardized way. This is done using data definitions. These definitions can be used across projects, and across different applications within the Cloud Pak for Automation. Each data definition has a title, description, and a datatype. We can also set a data definition as required or not. When a document is ingested into ADP, it results in a list of Key Value Pairs' (KVP) for that document. The Designer maps some of these KVP's to fields and teaches the model on how to extract the fields from the full list of KVP's. The designer then maps some of those fields to data definitions for a particular document type. Only the fields that have been mapped to data definitions will become Content Process Engine properties.

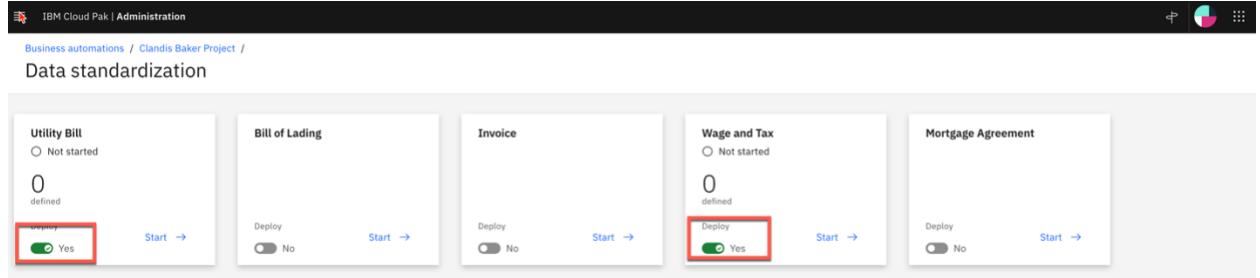
- _1. Return to the guided configuration flow and **Click** anywhere in the **Data standardization** box

Model Type	Status	Types Trained	Accuracy
Classification model	Ready	3	100%
Extraction model	Ready	3	97%
Data standardization	Not ready		

Here, you will see a list of available document types. Only the ones which have Deployed turned on will be visible in the verify interface and will have fields stored in FileNet.



_2. Ensure the Utility Bill and Wages and Tips and Deploy is toggled to **Yes**

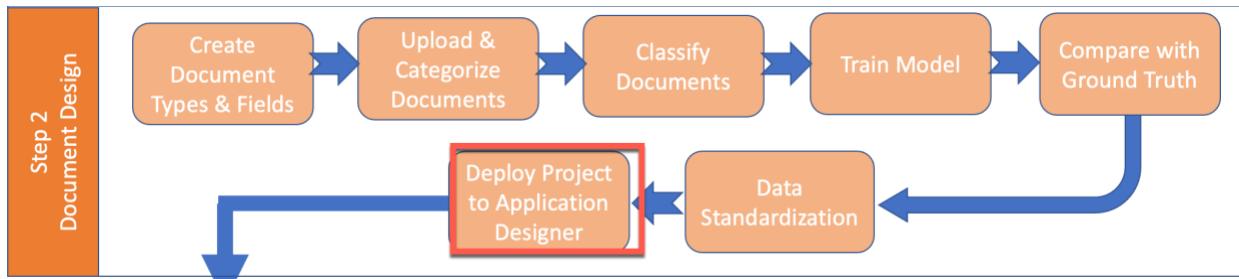


_3. Click on **Start** on either selected deployment.

This is where we begin defining the data file attribute definitions. You could create a new data definition and configure them. We will NOT be creating/defining any data fields for this lab.

_4. Return to the guided configuration screen by **Clicking on <your project>** name at the top of the screen.

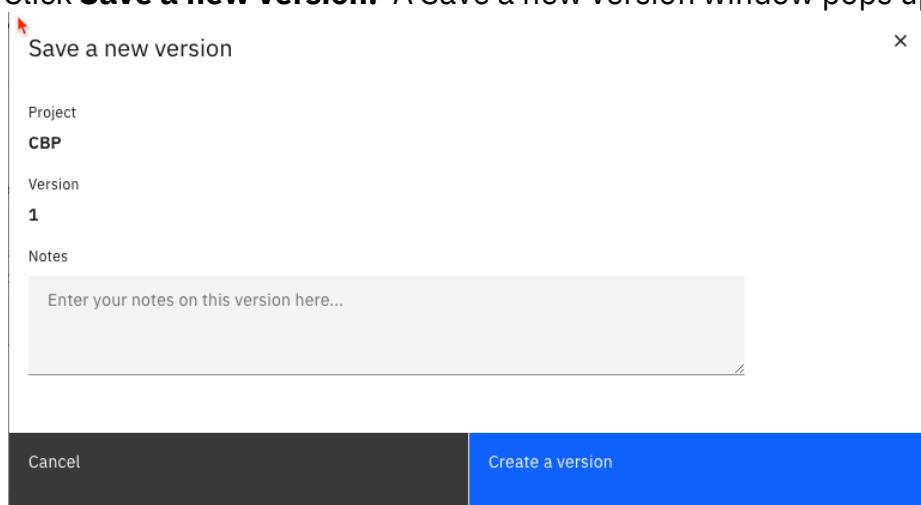
10 Version and deploy your project



At this point in our Designer project, we have defined a document type, labeled the fields we want from the document, trained (classified) the system to recognize the document type, reviewed the extracted fields we wanted and standardized (mapped) the document fields to our output.

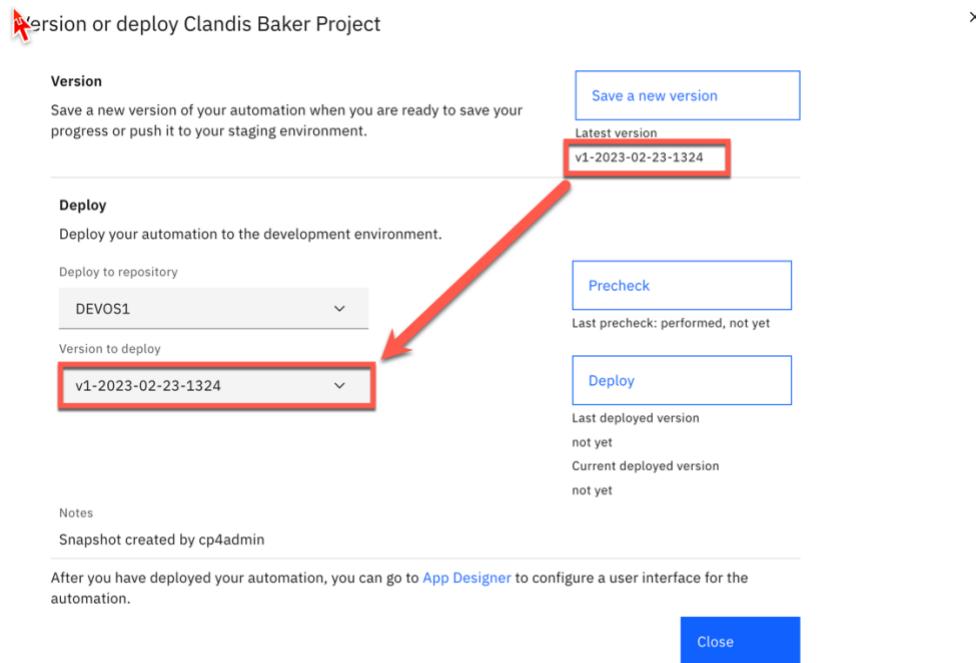
Now that we completed the configuration of the content extraction project, we need to save and deploy the design project to the application side. This will allow you to test your project using a client runtime interface.

- _1. If not already there, return to the guided home screen by clicking on your project name. Then **Click Version / Deploy**
- _2. Click **Save a new version**. A Save a new version window pops up.



- _3. Click on **Create a version**.

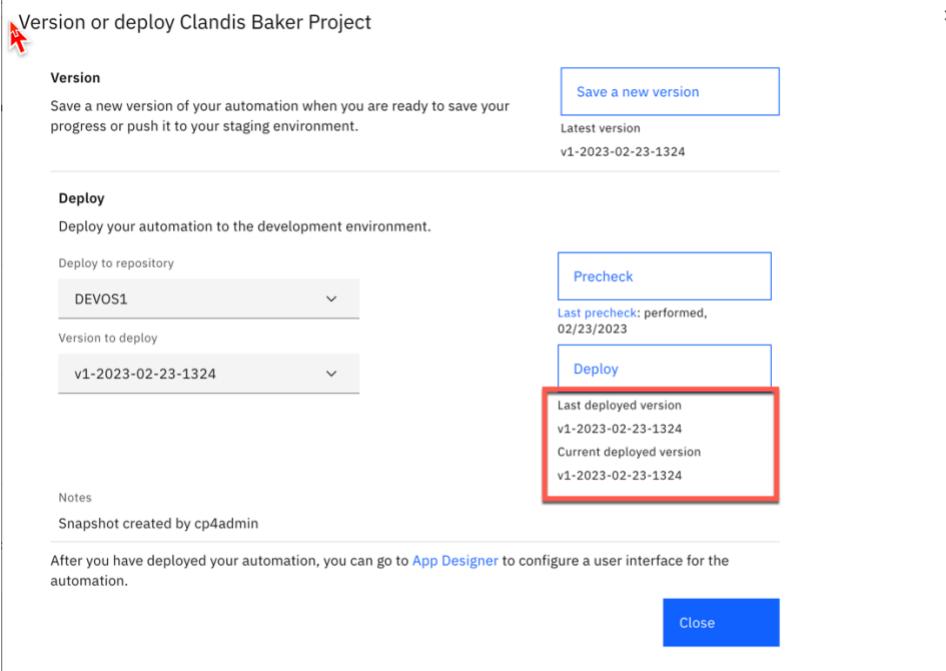
_4. Once the version is saved, you should see the version in the Version to deploy drop down list



... also, in the top corner has the “Latest Version.”

_5. Click on the Deploy button. This will also take a minute or two to deploy.

Once completed, you should have a notice that the project was deployed.





Note that you do not have to remain in the deploy screen while it is versioning or deploying. You can always click the button and then go back into any other screen if you like. It will run in the background. If you do this, just keep an eye on the top right of your screen for deployment status.

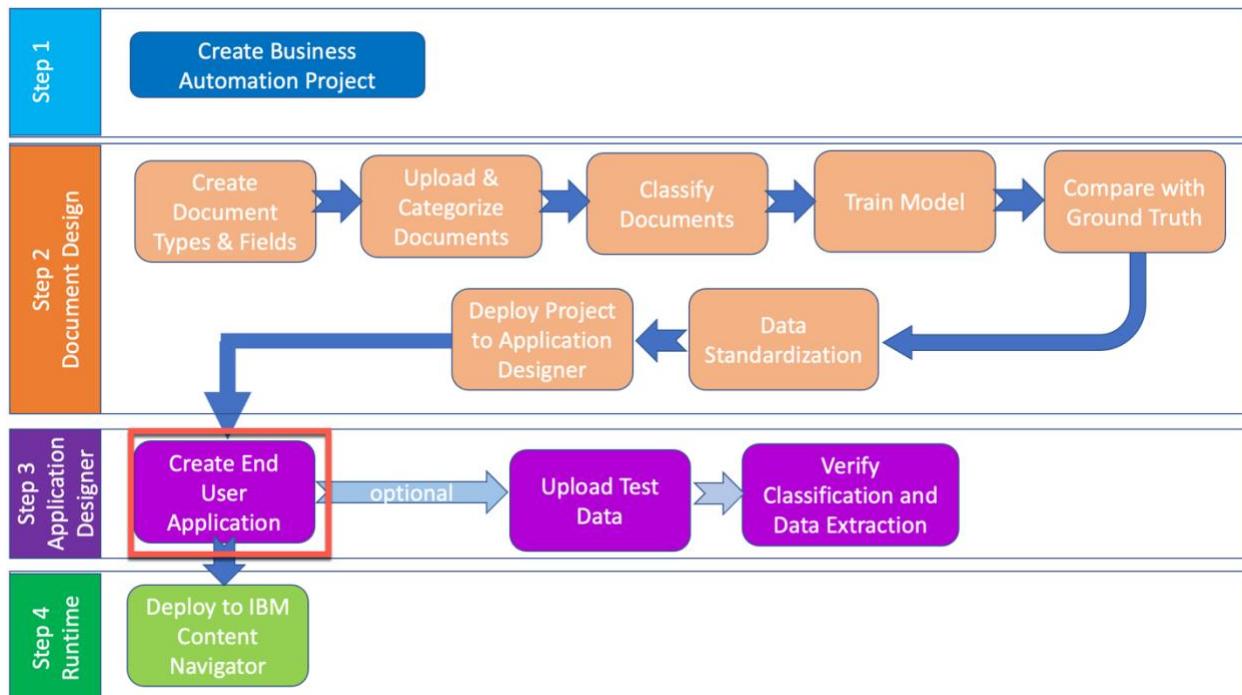
6. Click Close button.

Once deployed, proceed to the next steps.

From the home screen you can see the latest version and deployment

Document types and samples	Ready	5	20
Upload sample documents to define the types of documents you want the system to process.		types	samples on average
Classification model	Ready	5	100% accuracy
Train the model to classify your documents.		types trained	
Extraction model	Ready	4	
Train the model to extract the data from your documents.		types trained	
Data standardization	Not ready	0	
Map Fields to new or existing data definitions.		types reviewed	
Document retention	Ready	5	
Determine how long you want documents to stay in your content repository.		types reviewed	

11 Application designer



At this point we have designed or built a project that consists of document types, data or file types and methods to extract the desired data. The next major section of this lab is to build the user interface using the Application Designer. IBM provides two application templates for Document Processing

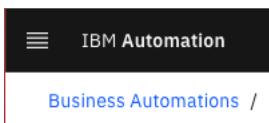
1. Batch Document Processing template – used to process batches of documents.
2. Document Processing Template – used to process single documents.

The lab will have you create a new batch processing application. We will quickly explore the various tabs in the interface, preview what the IBM Content Navigator (ICN) client would look like using the Preview feature and then publish our application to ICN where we will process a batch of documents.

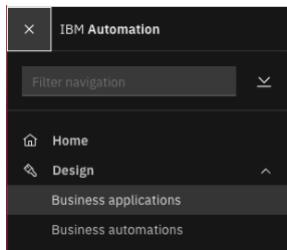
Changes to the application itself will not be in the scope of this lab.

11.1 Create your Runtime Application.

- _1. Return to the starting screen by **clicking the hamburger** in the top left.



and selecting **Business Applications**



_2. From the **Create** drop down list, select Application

The screenshot shows the 'Business applications' section of the IBM Cloud Pak Administration interface. The 'Create' dropdown menu is open, showing options: Application (selected), Template, and Toolkit. To the right, there is a message: 'There are no applications to display...yet. Click Create to get started. You can build an application by starting with a template and using shared assets in a toolkit.' Below this message, three template cards are displayed:

- Request Approval template**: Use this template to create a service desk request.
- Onboarding Application template**: Use this template to onboard new employees to your organization.
- Exception Handling template**: Use this template to create a basic refund request application.

_3. Select Enter your <application name> in the Name field.

The screenshot shows the 'Create a business application' dialog box. The 'Name' field is populated with 'Clandis Baker Application' and has a red border around it. The 'Purpose (optional)' section contains a placeholder text: 'Describe the purpose of the application'. The 'Create from template (optional)' section lists several templates under 'Select a base template': Exception Handling template (EHT), Onboarding Application template (OAT), Request Approval template (RAT), Document Processing template (CAT), and Batch Document Processing template (BCAT). At the bottom right of the dialog are 'Cancel' and 'Create' buttons, with 'Create' being highlighted in blue.

- _ 4. In the Create Form Template in drop down **select Batch Document Processing template (BCAT)**.

Create a business application

Name
Clandis Baker Application

Purpose (optional)
Describe the purpose of the application

Create from template (optional)

Select a base template

- Exception Handling template (EHT)
- Onboarding Application template (DAT)
- Request Approval template (RAT)
- Document Processing template (CAT)
- Batch Document Processing template (BCAT)**

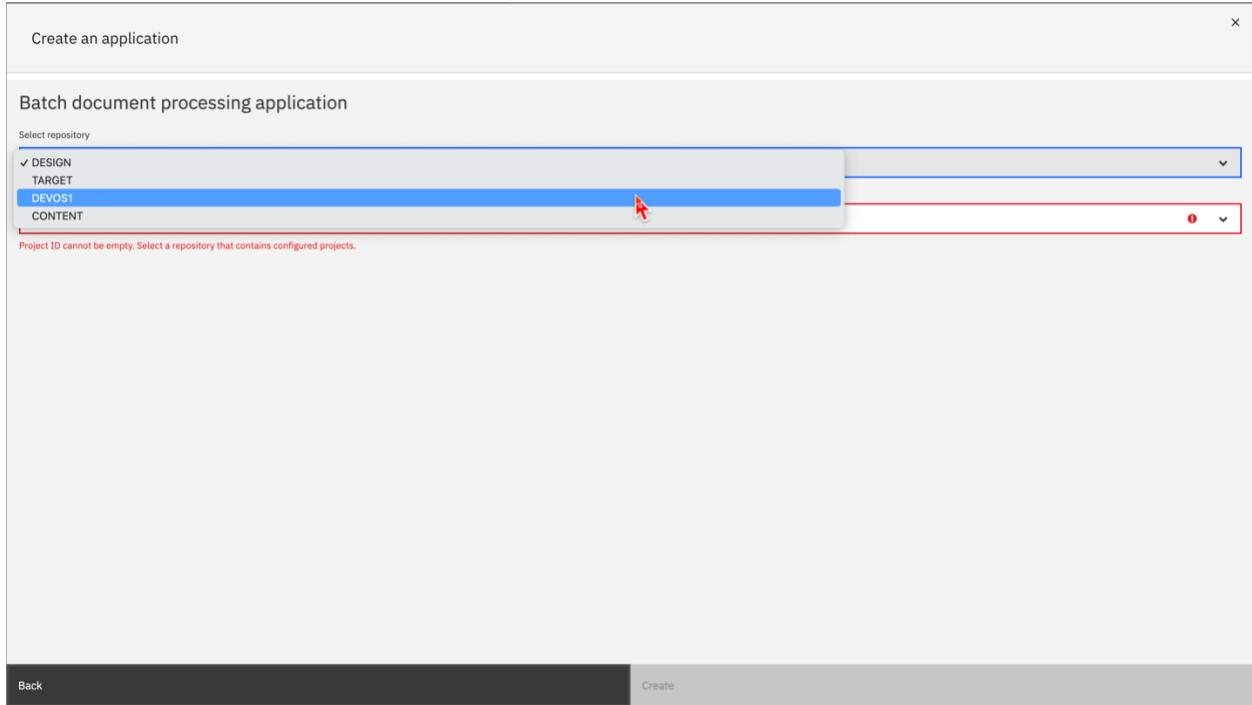
Cancel Create



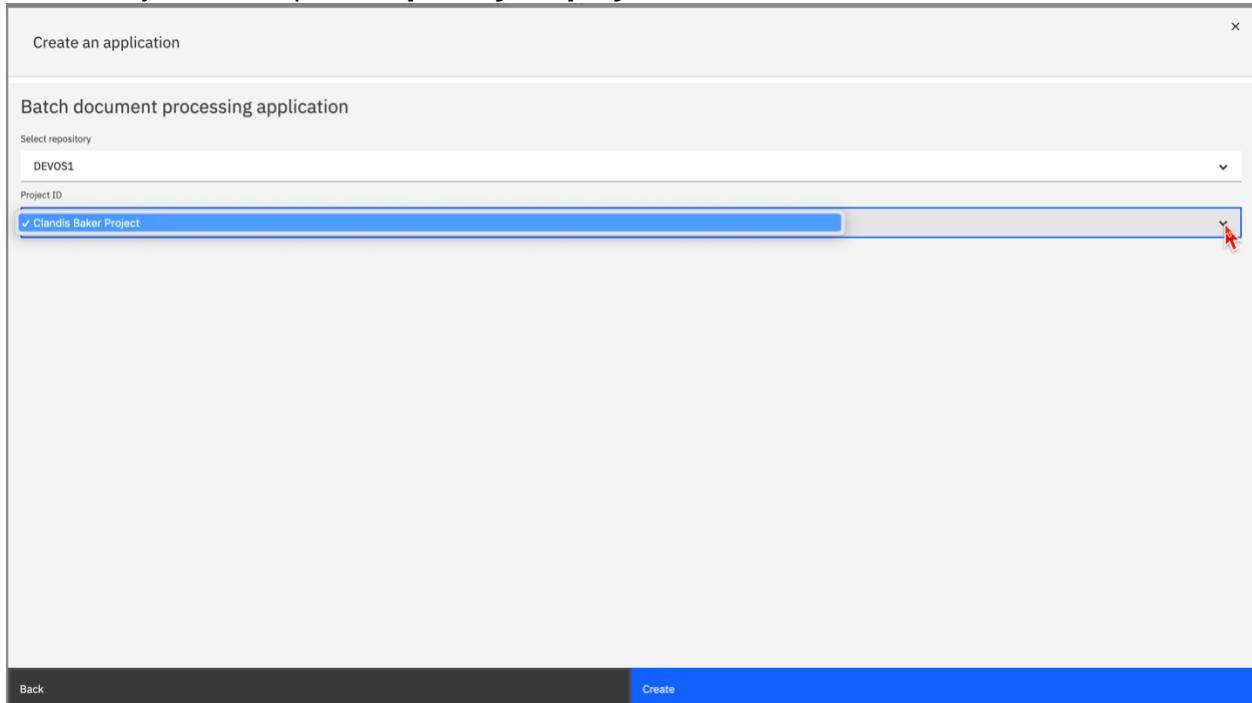
You could have selected the Document Processing Template if you only wanted to process a single document at a time, but in this lab, you will process several documents in a batch.

- _ 5. Click **Next**

- _ 6. You will be presented with the Create an application window. In the Select repository **pick DEVOS1**



_7. In the Project ID drop down **pick <your project name>**.



Note it may take a minute or two before this update and you can see your project

8. Click **Create**

You should now be in the *Application Designer*

The screenshot shows the IBM Cloud Pak | Administration interface for the 'Clandis Baker Application'. The top navigation bar includes 'IBM Cloud Pak | Administration', a user icon, and a 'Preview' button. Below the navigation is a breadcrumb path: 'Business applications / Clandis Baker Application'. The main content area displays a 'Review batch issues' section with two tabs: 'Document type and page order issues' and 'Data extraction issues'. To the right is a sidebar titled 'Drag a component to your page' containing a grid of icons for various components like 'Add batch modal', 'Content properties', and 'Document thumbnail'.

Name	Size	Modified by	Last modified	Version
My Document1	2 KB	User1	10/1/2022, 01:10 AM	1
My Document2	1 MB	User2	10/2/2022, 02:20 AM	2
My Document3	90 B	User3	10/3/2022, 03:30 AM	3



Batch Document Processing template (BCAT) has all the necessary pages and configuration to start using the application. Using this designer user interface, you have the option to further customize the application, such as its page design or actions, to fit your requirements.

9. Click on **Business applications** breadcrumb at the top.

The screenshot shows the IBM Cloud Pak Administration interface. The top navigation bar includes 'IBM Cloud Pak | Administration', 'Business applications / Clandis Baker Application', 'Preview' (with a note 'Last saved seconds ago by you.'), and a toolbar with icons for Content, Grid, and other functions. The main content area is titled 'Review batch issues' and contains two sections: 'Document type and page order issues' and 'Data extraction issues'. Below this is a 'Batches' section with a 'Content List' table:

Name	Size	Modified by	Last modified	Version
My Document1	2 KB	User1	10/1/2022, 01:10 AM	1
My Document2	1 MB	User2	10/2/2022, 02:20 AM	2
My Document3	90 B	User3	10/3/2022, 03:30 AM	3

Items per page: 100 Items 1-3

A sidebar on the right is titled 'Drag a component to your page' and lists various UI components with their icons and names, such as 'Add batch modal', 'Content list', 'Content properties', 'Delete object modal', 'Document correction', and 'Document thumbnail'.



Note: It may take several seconds to build and display the current configuration of the interface.

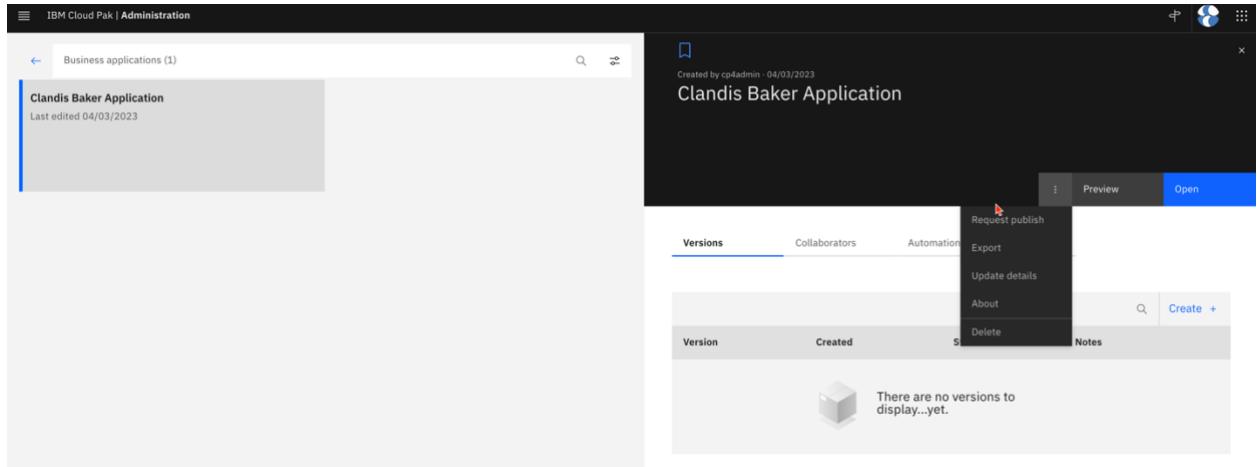
- _10. From this screen you can also select preview to see the pre-configured interface. The open takes you the application where you can modify the look and feel. If you hover anywhere in the box it will turn grey and **Click**.

The screenshot shows the 'Business applications' screen. The top navigation bar includes 'IBM Cloud Pak | Administration', 'Business applications (1)', and a search bar. The main content area shows a single application card:

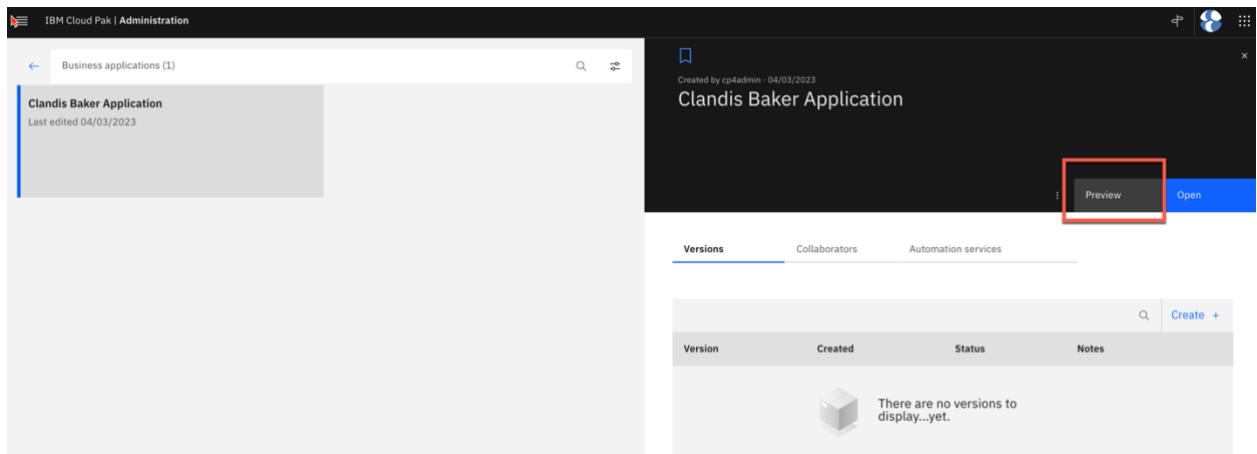
Clandis Baker Application
Last edited 04/03/2023

Below the card are 'Preview' and 'Open' buttons. The sidebar on the left has a title 'Business applications' and includes a 'Create' button, an 'Import' button, and links to 'Applications', 'Templates', and 'Toolkits'.

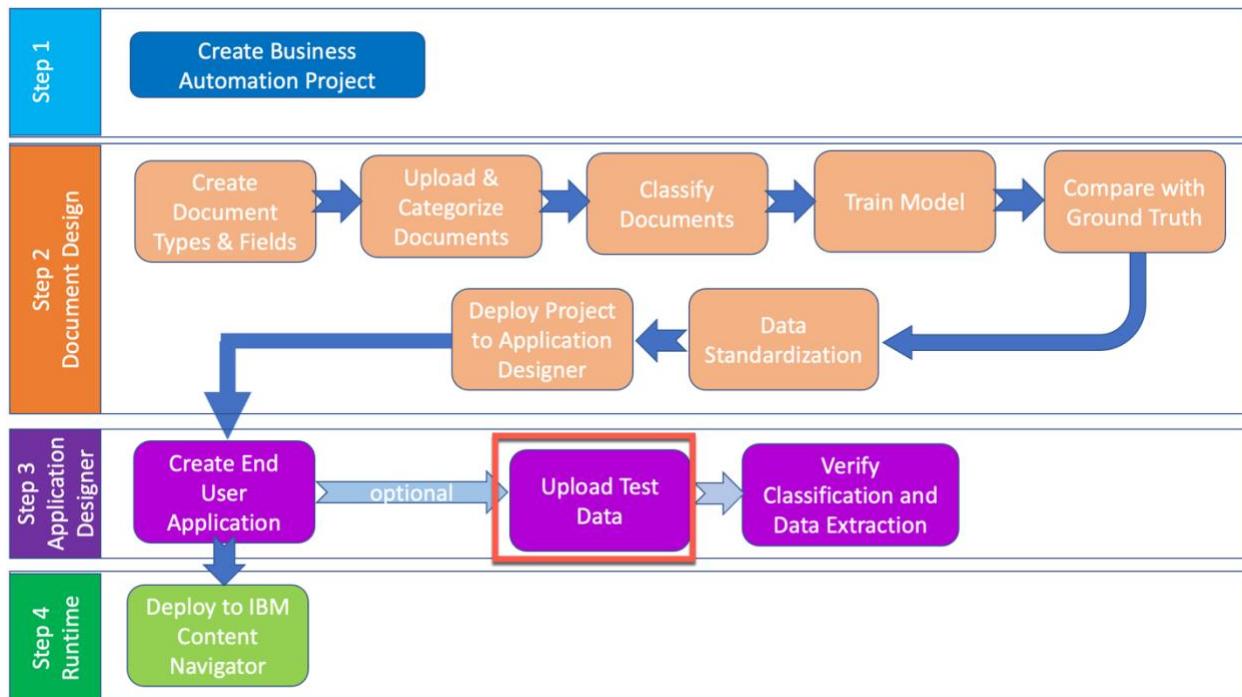
- _11. From this screen if you **click** on the **3 dots** you can save versions, export application, or delete the application. Just wanted to show this for future reference.



12. Click on Preview



11.2 Upload documents for processing



- _1. You should be in the default application user interface for ADP it opens a new tab in your browser.

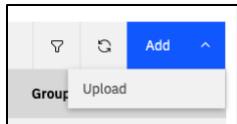
The screenshot shows the 'Review batch issues' section of the ADP application. It displays two categories of issues:

- Document type and page order issues:** 0 batches
- Data extraction issues:** 0 batches

Below this, there is a search bar and a message stating "No items found."

There are two key screens you will work with: “*Document type and page order issues*” and the “*Data extraction issues*”. First, we need to upload some test documents and have them processed.

- _2. Click on Add, then Upload.



- _3. Enter a **name** for your batch in the Display Name field and set the **Priority to High** as seen in the image below.

Upload new batch

* Display Name
Batch 1

Description

Priority
High

- _4. Click **Select files**.

Navigate to the samples folder previously downloaded from [Section 2](#) and use the *Group 3 - Runtime Demo Set* folder documents. Select all the files in the folder.

- _5. Click **Open**

You will see a window that will give the operator a chance to manually classify the documents before they are ingested. By clicking on one of the files you will be presented with an option to manually classify the documents. In the example below would be how to manually classify a document.

Add Files

To manually specify document type, first select the files in the table. Use the classify option, to assign the document type for selected file(s). If a file is not manually classified, the system will auto-classify it.

1 items selected		Classify ▾	Auto Classify	Deselect
<input type="checkbox"/>	File Name	Document Type		
<input checked="" type="checkbox"/>	B_PO_5.pdf	Auto Classify		
<input type="checkbox"/>	DE_FW2_1000_0001F.pdf	Auto Classify		
<input type="checkbox"/>	DE_FW2_4000_0011F.pdf	Auto Classify		
<input type="checkbox"/>	DE_FW2_4001_0001S.pdf	Auto Classify		
<input type="checkbox"/>	DE_FW2_4001_0010F.pdf	Auto Classify		

Cancel **Add**

We are not going to do this but instead let ADP auto classify them.

Add Files

To manually specify document type, first select the files in the table. Use the classify option, to assign the document type for selected file(s). If a file is not manually classified, the system will auto-classify it.

<input type="checkbox"/>	File Name	Document Type
<input type="checkbox"/>	B_PO_5.pdf	Auto Classify
<input type="checkbox"/>	DE_FW2_1000_0001F.pdf	Auto Classify
<input type="checkbox"/>	DE_FW2_4000_0011F.pdf	Auto Classify
<input type="checkbox"/>	DE_FW2_4001_0001S.pdf	Auto Classify
<input type="checkbox"/>	DE_FW2_4001_0010F.pdf	Auto Classify

**_6. Click on the Add button.**

A screenshot of the software interface showing a progress bar at the bottom of a table. The progress bar is partially filled, with the text '3 of 5 files processed' next to it. The table header includes columns for Name, Files, Priority, Status, Added on, Added by, Group, and Location. A blue 'Add' button is visible in the top right corner of the table area.

A progress bar will be displayed indicating when all documents have been uploaded.

_7. Click the 3 dots at the end of the line.

A screenshot of the software interface showing the same table as before. A red arrow points to the three-dot menu icon (three vertical dots) located at the far right end of the table's header row. The table data remains the same, showing one item named 'Batch01'.

_8. Click Submit

In the screen shot below, you see we have a document issues (status) and we now have 1 batch in the “Document type and page order issue” tile.

Name	Files	Priority	Status	Added on	Added by	Group	Location
Batch01	5	High	⚠ Document issues	03/27/2023, 01:45 PM	cp4admin		

11.3 Correct any classification errors.

1. Click on the Document type and page order issues tile to open the batch.

Name	Priority	Status	Added on	Added by	Group	Location
Batch 1	High	Document issues	01/13/2021, 08:44 am	CEAdmin		

2. Click on <your batch name> to open it.

You should now see all the documents you uploaded in your batch. The ones with issues will have a yellow checkmark for documents that have a low confidence document type and a red exclamation mark for documents it could not classify.

Batch01

Cancel Save changes Submit

Documents (5)	
Document name	Document type
B_PO_5.pdf	Undefined
BAD_FW2_1000_0003F.pdf	Wage and Tax
TR_FW2_1000_0003F.pdf	Wage and Tax
TR_FW2_2000_0003F.pdf	Wage and Tax
TR_FW2_4000_0002F.pdf	Wage and Tax



PURCHASE ORDER
RUBE'S Meat Co.
P.O. No.: 71238
DATE: 09 March 2020
CUSTOMER ID: 46273CD

VENDOR:
Chicken Run Ranch
24 Quay Street
Nelson Village NE23 8DD
UK
078-2054-8486

SHIP TO:
Rube's Meat Co.
64 Pendlevalt Road
Burton, Leonard HG3 2SU
UK
078-7875-2017

SHIPPING METHOD	SHIPPING TERMS	DELIVERY DATE			
AIR	C.I.F.	29 March 2020			
QTY	ITEM #	DESCRIPTION	JOB	UNIT PRICE	LINE TOTAL
250 PCS	01	Whole Chicken		£1.50	£345.00
150 Packs	02	One Day Old Chick		£1.00	£157.50

TOTAL: £502.50

- _3. Most of the document types are correct but it looks like a PO got mixed into our batch so we can **Click on the Trash can** to delete it from the batch. And **select OK** to delete it.

Batch01

Cancel Save changes Submit

Documents (5)	
Document name	Document type
B_PO_5.pdf	Undefined
BAD_FW2_1000_0003F.pdf	Wage and Tax
TR_FW2_1000_0003F.pdf	Wage and Tax
TR_FW2_2000_0003F.pdf	Wage and Tax
TR_FW2_4000_0002F.pdf	Wage and Tax



PURCHASE ORDER
RUBE'S Meat Co.
P.O. No.: 71238
DATE: 09 March 2020
CUSTOMER ID: 46273CD

VENDOR:
Chicken Run Ranch
24 Quay Street
Nelson Village NE23 8DD
UK
078-2054-8486

SHIP TO:
Rube's Meat Co.
64 Pendlevalt Road
Burton, Leonard HG3 2SU
UK
078-7875-2017

SHIPPING METHOD	SHIPPING TERMS	DELIVERY DATE			
AIR	C.I.F.	29 March 2020			
QTY	ITEM #	DESCRIPTION	JOB	UNIT PRICE	LINE TOTAL
250 PCS	01	Whole Chicken		£1.50	£345.00
150 Packs	02	One Day Old Chick		£1.00	£157.50

TOTAL: £502.50

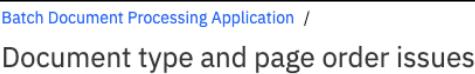
- _ 4. Review all documents to ensure everything is correct. If the system no longer detects any issues, you should see a green checkmark near the top of the document list.



- _ 5. Click **Save Changes** and then **Submit** to save your changes and have the batch processed.

The system will start reprocessing the documents now that they have been classified correctly.

- _ 6. Click on the blue **Batch Document Processing Application link** at the top to return to the previous preview menu.



11.4 Correct extraction issues

The following instructions are based on a pre-trained sample application. Not what you will see in your untrained application.



Important Note: The project you are using for this has been configured but NOT run through the training (Deep Learning). So, the results will not reflect what they should be. IN A NORMAL SCENARIO, ON A CLUSTER WITH GPU AND DEEP LEARNING ENABLED, YOU WOULD HAVE TRAINED YOUR MODEL BEFORE DEPLOYING IT AND WOULD BENEFIT FROM HIGHER EXTRACTION RATES. the purpose of this lab is to teach you the tools but won't show you the trained results.

It may take a few seconds for your batch to advance to the next step. If your batch needs further attention, you will see it appear in the Data extraction issues tile.

- _ 1. Click on the **Data extraction issues** tile to open it.



_2. Click on <your Batch name> to open.



After opening we see all the documents that have been processed but one looks to have extraction issues.

Name	Issues	Status	Modified on	Modified by
BAD_FW2_1000_0003F.pdf	1	Data issues	03/04/2023	cp4admin
TR_FW2_1000_0003F.pdf		Issues reviewed	03/04/2023	cp4admin
TR_FW2_2000_0003F.pdf		Issues reviewed	03/04/2023	cp4admin
TR_FW2_4000_0002F.pdf		Issues reviewed	03/04/2023	cp4admin

Items per page: 100 1-4 of 4 items

_3. Click on the bad document to open it. Zoom in a bit to get a better picture of the document.

Extracted data

All Fields

Federal Income Tax Withheld

Federal Income Tax Withheld
9000.00

Employee Social Security Number *

(none)

Employer Identification Number

(none)

Employers Name and Address

Bricks and Mortar 343 Jackson Ave Costa Mesa, CA 90394

Social Security Wages

75000.00

Wages Tips Other Compensation

(none)

Employee Name and Address

Last name Suff. Stella K. James 343 Twisting Way Red Beach, CA 90354 f Employee's address and ZIP code

Take a moment to discover the image viewer features.

Image viewer features at top:

The screenshot shows a document processing interface. On the left is a thumbnail view of the W-2 form. The main area displays the W-2 form with various fields highlighted by a red box, including the employee's name, address, and tax withheld amounts. To the right is a panel titled "Extracted data" which lists the extracted fields and their values. A validation warning is shown for the Employee Name and Address field.

Field	Value
Employee Name & Address	Francis A. Hellbut 123 Main Street New York, NY 10028
Federal Income Tax Withheld	44025.44
Social Security Wages	132099.00
Medicare Wages and Tips	2369.87
Total Wages, Etc.	163439.33
State Income Tax Withheld	3092904
Total Income Tax Withheld	183439.33
Local Income Tax Withheld	521.00
Locality Name	Red Beach

- Rotate image.
- Visual effect adjustment
- Invert

Image viewer features at bottom:

This screenshot shows a different view of the same W-2 form. The main area displays the W-2 form with various fields highlighted by a red box, including the employee's name, address, and tax withheld amounts. To the right is a panel titled "Extracted data" which lists the extracted fields and their values. A validation error is shown for the Employee Social Security Number field.

Field	Value
Employee Name & Address	Stella K. James 343 Twisting Way Red Beach, CA 90354
Federal Income Tax Withheld	9000.00
Social Security Wages	75000.00
Medicare Wages and Tips	1200.00
Total Wages, Etc.	75000.00
State Income Tax Withheld	2250.00
Total Income Tax Withheld	75000.00
Local Income Tax Withheld	45.00
Locality Name	Red Beach

- Page and thumbnail's view

- Fit to window.
- Zoom and Magnify

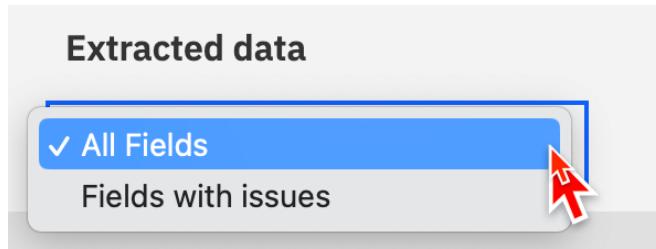
Field features

The screenshot shows a document processing application interface. On the left, there is a preview of an W-2 form. The right side displays the extracted data from the form. A dropdown menu labeled "Extracted data" is open, with "All Fields" selected. Below this, the "Federal Income Tax Withheld" field is shown with the value "9000.00". A validation error icon (a red circle with a exclamation mark) is present next to the "Employee Social Security Number" field, which has "(none)" listed. Other fields like "Employer Identification Number", "Employers Name and Address", "Social Security Wages", and "Wages Tips Other Compensation" are also listed with their respective values or "(none)".

- Show all fields.
- Show fields with issues.

Also note that fields that have issues have a notification icon next to them. For example, Wages Tips Other Compensation field picked up correctly but has a low confidence based on the extraction results.

4. Under Extracted data click on the drop down twisty.



5. Click on the ALL Fields.

This view shows all the fields that we defined earlier. Fields with an asterisk are mandatory fields.

Change the Extracted data back to Fields with issues.



The Employee Social Security Number is a mandatory field. For purposes of this lab it was changed to “Bad SSN”. Since you did not make that phrase an alias ADP was not able to pick it up.

_6. Click on Employee Social Security Number and with your mouse select the SSN under “Bad SSN”.

The screenshot shows the ADP interface with a W-2 form. A red box highlights the 'Employee Social Security Number' field, which contains the value '183-94-7103'. To the right of the form, a sidebar titled 'Extracted data' shows a list of fields with issues. The 'Employee Social Security Number' field is listed with the message 'Required value is missing.'

Also the Wages Tips Other Compensation did not have a correct alias defined. But since it was not a required field, you can continue to process.

_7.. Click on Save Changes box at the top.

The screenshot shows the top navigation bar of the ADP interface. The 'Save changes' button is highlighted with a blue border. Other buttons include 'Cancel', 'Done and next', and 'Done'.

_8. For the remaining fields there are no extraction issues that ADP picked up for mandatory fields. You may see some low confidence characters. If so, Click on Dismiss for each field with a yellow validation warning.

_9. Click on Done and next.

_10. All documents have been processed **Click** on **Submit** at the top to complete the batch.

12 Optional Export/Import Project.

If you would like to save your project and perhaps use it later, you can do this lab.
From the Business Automations

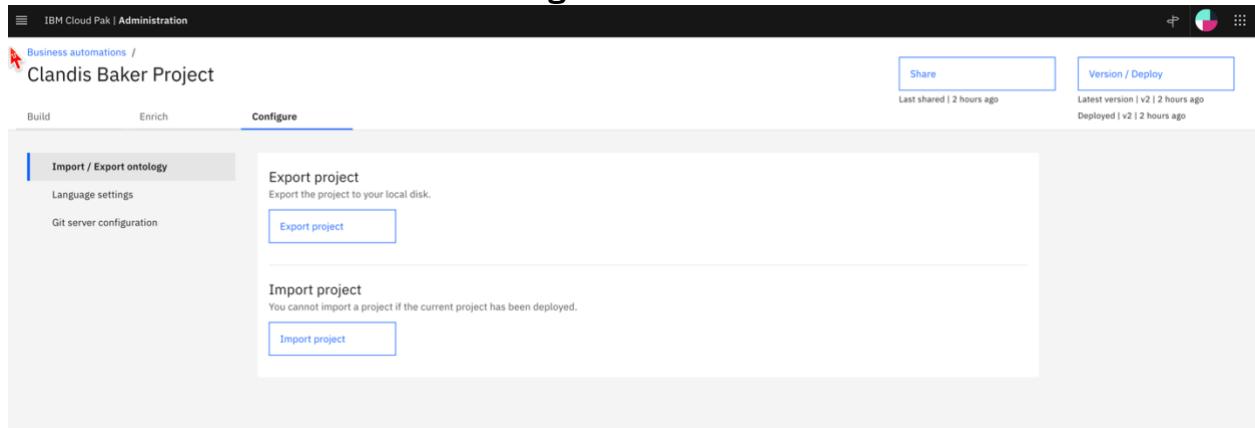
_1. From the Business Automations screen **select Document Processing**.

The screenshot shows the 'Business automations' screen in the IBM Cloud Pak interface. At the top, there's a navigation bar with 'IBM Cloud Pak | Administration'. Below it, a sidebar on the left lists categories: 'Published automation services' (with 'Decision', 'Document processing', 'Workflow', and 'External' options), 'Decision', and 'External'. The 'Document processing' option is highlighted with a blue border. The main content area on the right shows a single item: 'Clandis Baker Project' (Last edited 02/23/2023). At the bottom, there are 'Create', 'Import', and download buttons.

_2. Select <your project name> Click open

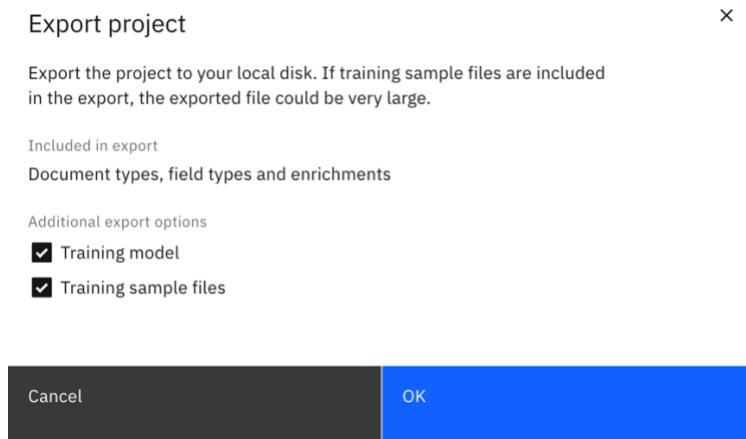
This screenshot is from the same interface as the previous one, showing the 'Clandis Baker Project' card. A red hand icon with the word 'Open' is overlaid on the 'Open' button in the top right corner of the project card.

_3. From the Main screen select the Configure tab



_4. Select Export Project

_5. On Export Project window check Training Module and Training Sample files



_6. Click on OK

_7. A project-export-<date-time>.zip will be download via browser to local machine.

END OF LABS

Appendix A - Troubleshooting

TechZone Pending Status taking Long Time

Operator shows Pending status in a namespace – OLM know issue.

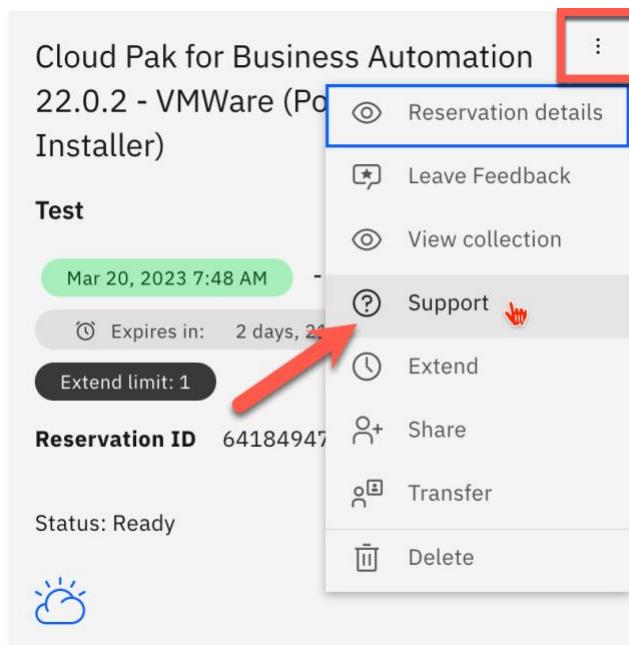
An operator fails to install and continuously shows Pending status.

For fix visit below link.

<https://www.ibm.com/docs/en/cpfs?topic=ii-operator-shows-pending-status-in-namespace-olm-known-issue>

Other issue could be the deployment itself had an issue. Two things to do in this case.

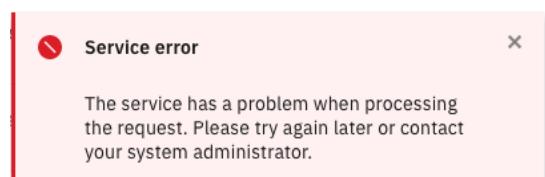
1. Open a support ticket by clicking on the 3 dots on the tile.



IBM Internal can also access support via SLAC Channel at #itz-techzone-support

2. Delete tile and try to deploy again.

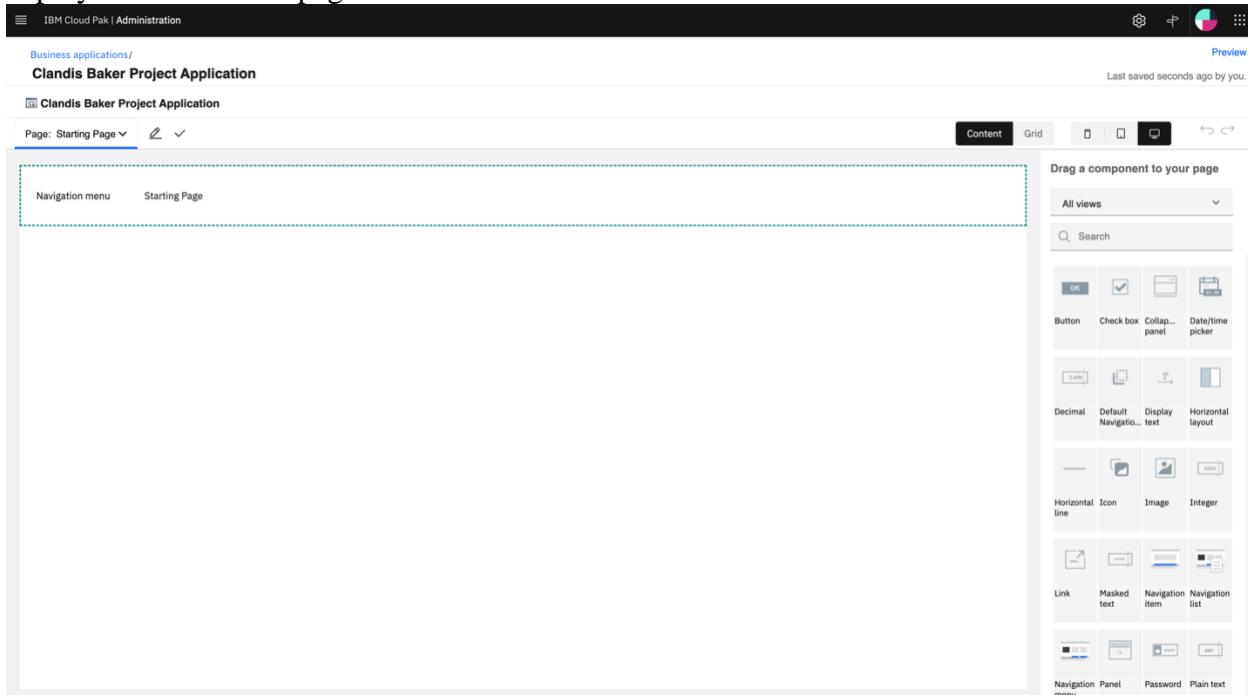
Service Error



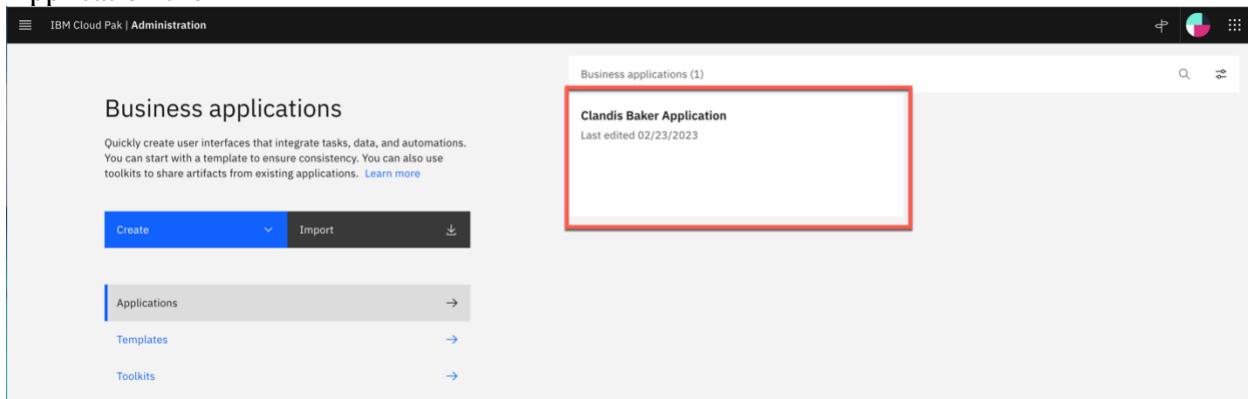
There was trouble communicating most of the time you can simply ignore and continue on.

Application Blank

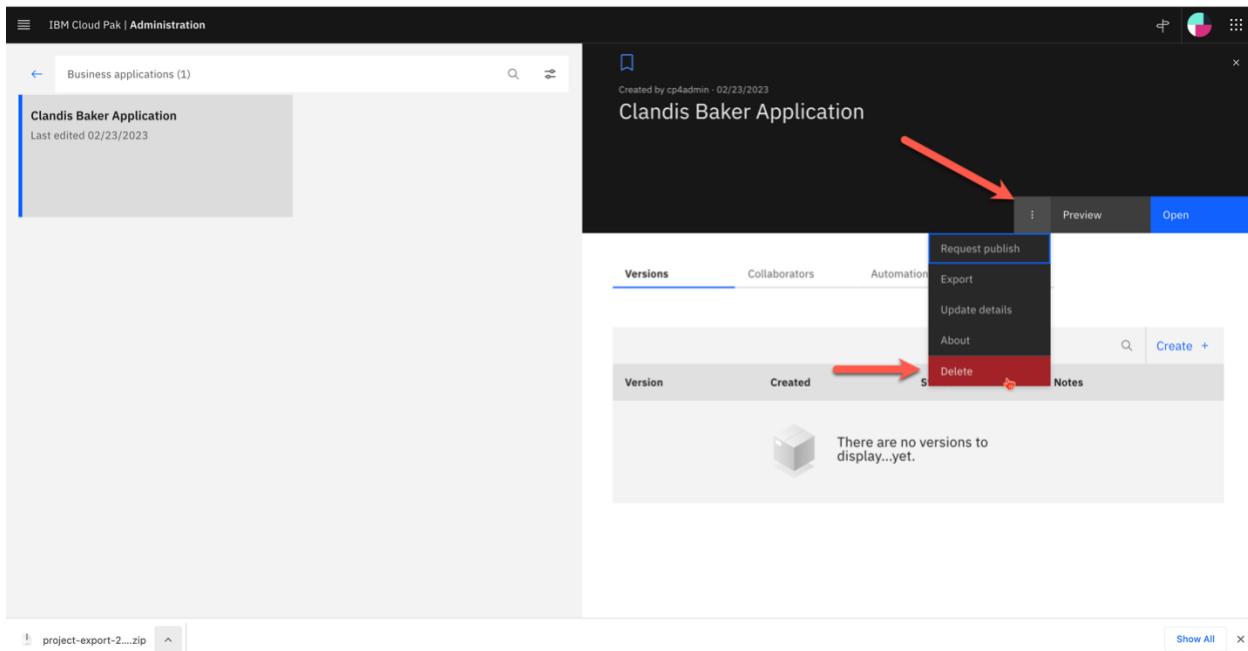
During creating of Business Application setup, sometimes on first time after project has been deployed. The Starter page is blank.



If this happens delete the application and try again. To delete the application, Click on the Application tile



Then Click on the 3 dots and Select Delete



Connection issue with Workstation to Cloud.

If issues with connection from workstation to cloud after it's been working. Reboot your workstation.

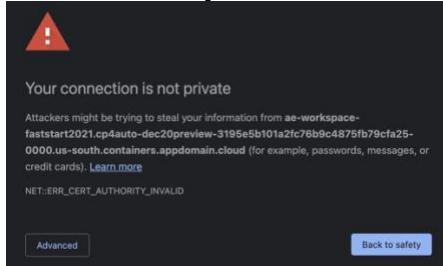
Opening an Incognito Window

When you open a new incognito window, you will need to accept certificates before logging in to ADP. Customers shouldn't have this issue because they will have their own certificates instead of the self-signed certificates used in this environment.

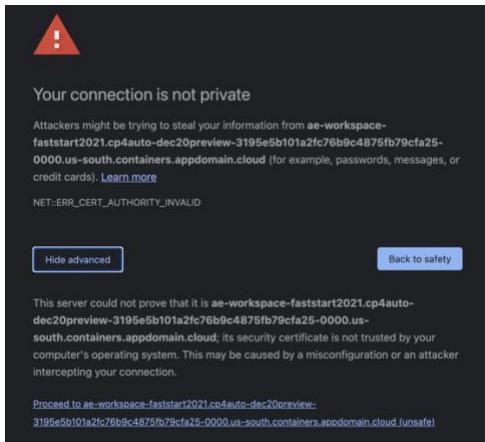
In your incognito window, go to the following URLs located in this Box:

Open the Generate Security Tokens Box note and click all 3 of the links listed. This will reset the self-signed security certificates.

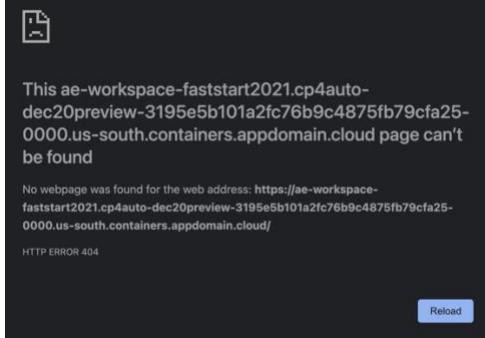
For each URL, your browser window will show a message like this:



Click Advanced, and the browser window will look something like this:



Click the “Proceed to...” link. You’ll see a message like this in your browser window:

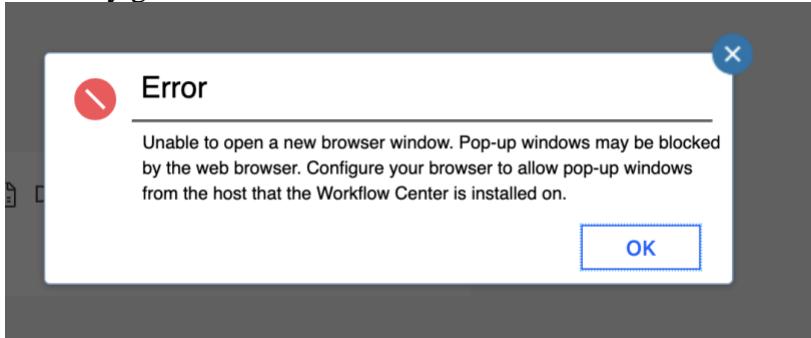


Ignore the error and proceed to the next link.

After doing this for each of the URLs above, log in to BAStudio

Popup Blocked when trying to Preview Application.

You may get error like this:



You will need to grant access to pop up windows in your browser.

Appendix B - BAW & ADP Integration Sample

For the End to End demo BAW was integrated with ADP. This link explains how to accomplish.
<https://github.com/IBM/baw-adp-integration-sample>

Appendix C - Badge Information.

Badge quiz page - <https://learn.ibm.com/user/policy.php>