

IBM Cloud Pak for Business Automation Demos and Labs 2024

Lab Guide – Automation Document Processing

V 2.2.8 (for CP4BA 23.0.2 IF002)

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)
krishtkrish@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	5
1.1 Icons.....	5
1.2 Abstract	5
1.3 Introduction.....	5
2 Getting started.....	7
2.1 Reserving an environment.....	7
2.2 Open your IBM Cloud Environment.....	10
2.1 Preparing your Jam-in-a-Box environment	12
3 Lab Overview	16
3.1 How does ADP work?	16
4 Create Document Processing Project.....	19
4.1 Reviewing the interface.....	24
4.1.1 Build Tab	25
4.1.2 Enrich Tab	25
4.1.3 Configure Tab.....	26
5 Configure a Wage and Tax document type	29
5.1 Create Wage and Tax document type	29
5.2 Create Field	31
5.3 Create the Employee Name Address field	36
5.4 Create Employee Social Security Number Field	37
6 Document Types and Samples Overview	40
6.1 Categorize documents.....	41
7 Train classification	47
7.1 How do I improve my results?.....	52
7.1.1 Option 1 – Add more samples	52
7.1.2 Option 2 – Review all uploaded samples	53
8 Data extraction	54
8.1 Correcting extracted values	57
8.2 Train extraction model	62
9 Data standardization	63
10 Version and deploy your project.....	65
11 Application designer	68
11.1 Create your Runtime Application.....	68
11.2 Upload documents for processing.....	75
11.3 Correct any classification errors	78
11.4 Correct extraction issues	80
12 Export/Import Project (Optional)	85
Appendix A - Troubleshooting.....	87
Application Blank	87
Popup Blocked when trying to Preview Application	88
Appendix B - BAW & ADP Integration Sample	89

1 Overview

1.1 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.

1.2 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.

1.3 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. It is intended as an entry point.

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/23.0.2/Document%20Processing/Lab%20Data>

- _1. Click on “Group1 – Design Docs for Tax Lab.zip”.
- _2. Then Click on Download

Name	Last commit message	Last commit date
...		
Group 1 - Design Docs for Tax Lab.zip	lab data folders reorg	6 minutes ago
Group 2 - Classification Results Increase Set.zip	lab data folders reorg	6 minutes ago
Group 3 - Runtime demo Set.zip	lab data folders reorg	6 minutes ago

- _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.

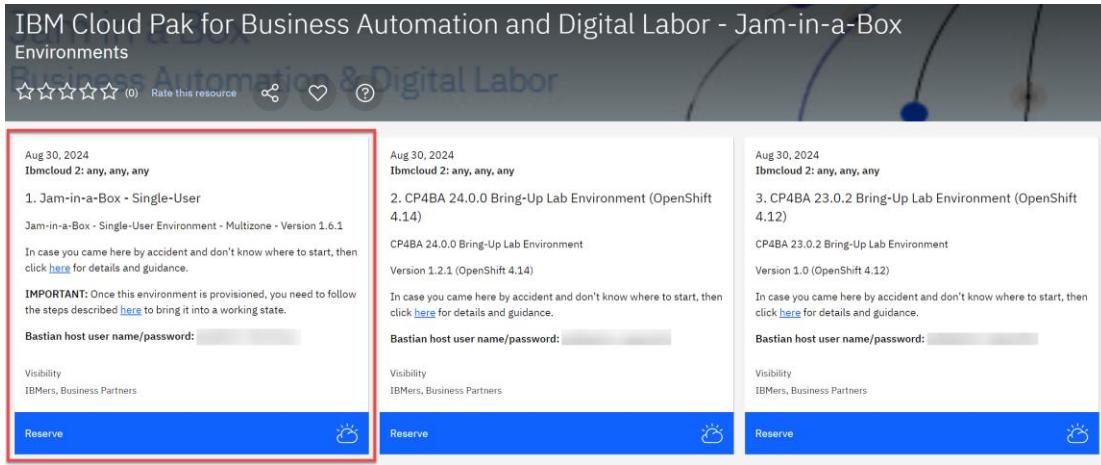
You will notice the images are in various unique folders that will be referenced specifically in the different labs later. Please keep them in their proper folders.

2.1 Reserving an environment

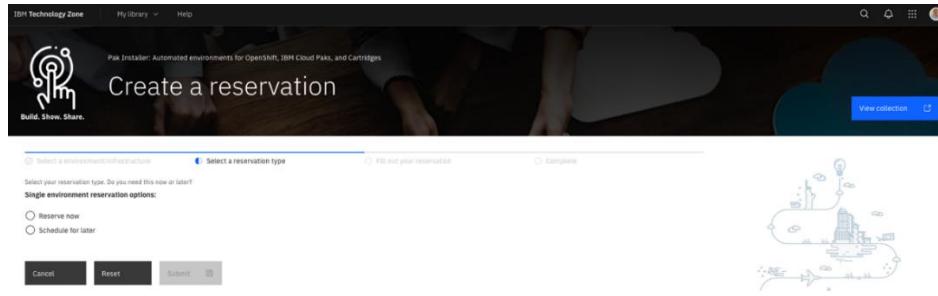
Below you'll find a description of how to reserve a Jam-in-a-Box environment from IBM TechZone. In case you don't have access to TechZone or have your own CP4BA environment with the appropriate version and ADP installed, you should be able to also perform the lab. Depending on the environment the name of the object store may differ.

- _1. Navigate to [IBM Cloud Pak for Business Automation and Digital Labor - Jam-in-a-Box](#)

Note the user name and password provided at the bottom of the **1. Jam-in-a-Box - Single-User** tile. These are later required to logon to the bastian host when you connect to it via RDP as described in section 2.2.



- _2. Click Reserve for the 1. Jam-in-a-Box – Single-User tile
- _3. On Create a reservation – Select a reservation type screen select option for when to start provisioning



- _4. On Create a reservation – Fill out your reservation screen select a purpose.

The selected purpose defines how long the environment is available, how many times it can be extended and if you need a Sales Opportunity number or not.

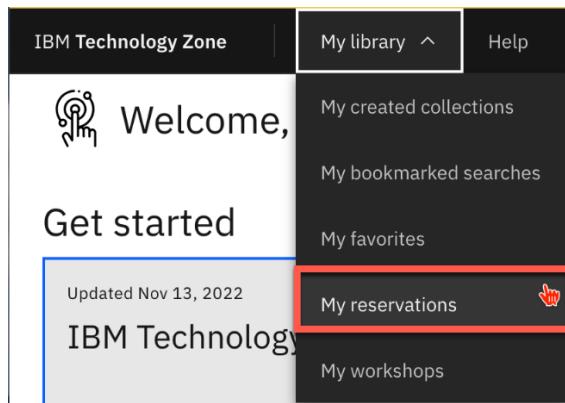
- _5. Enter <some description> in the Purpose description box.
- _6. For Preferred Geography (required) select your preferred data center location

- _7. Do enable the VPN access option

- _8. On the right hand pane, click on the option “I agree” and click on the “Submit” button.

The screenshot shows the 'Create a reservation' page. At the top, it says 'IBM Technology Zone' and 'My library'. The main area has tabs for 'View collection' and 'Edit'. Below that, there are sections for 'Name' (Jam-in-a-Box - Single-User) and 'Purpose'. Under 'Purpose', there are three options: 'Demo' (selected), 'Education' (Gaining experience with specific technology, product, or solution), and 'Test' (Need to test a specific function, configuration, or customization). There's also a note about selecting the correct purpose. Below 'Purpose' is a 'Sales Opportunity number' field and a note about providing an opportunity number. At the bottom, there's a 'Purpose description' field containing 'To learn ADP'. A checkbox for 'I Agree to IBM Technology Zone's Terms & Conditions and End User Security Policies' is checked, and a 'Submit' button is visible.

- _9. After about 1.5 to 2 hours, you should receive an email that your environment is ready. This is proceeded by an email saying that the environment is provisioning.
- _10. 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**, or by clicking the link in the email.



You can also access directly using the link below

<https://techzone.ibm.com/my/reservations>

2.2 Open your IBM Cloud Environment

- _1. Once the environment is created, you can open the reservation and you shall find the screen below

IBM Technology Zone | My library | Help

My reservations / Collection

Jam-in-a-Box - Single-User
To learn ADP

Date: Aug 22, 2024 8:53 AM | Aug 24, 2024 10:10 AM | Expires in: 1 days, 23 hours, 52 minutes | Extend limit: 2

Status: Ready

Published services

Bastion Host Remote Desktop (copy value to RDP): <useast.services.cloud.techzone.ibm.com:20370>

Purpose

Purpose	Opportunity ID(s)
Education	
Opportunity Product(s)	Opportunity description
Customer(s)	To learn ADP

Environment

Reservation ID: [Tuna](#)

- _2. Here you would notice the remote desktop url to the bastion host. Use the Microsoft remote desktop feature or a Remote Desktop client of your choice to connect to the highlighted URL.

IBM Technology Zone | My library | Help

My reservations / Collection

Jam-in-a-Box - Single-User
To learn ADP

Date: Aug 22, 2024 8:53 AM | Aug 24, 2024 10:10 AM | Expires in: 1 days, 23 hours, 51 minutes | Extend limit: 2

Status: Ready

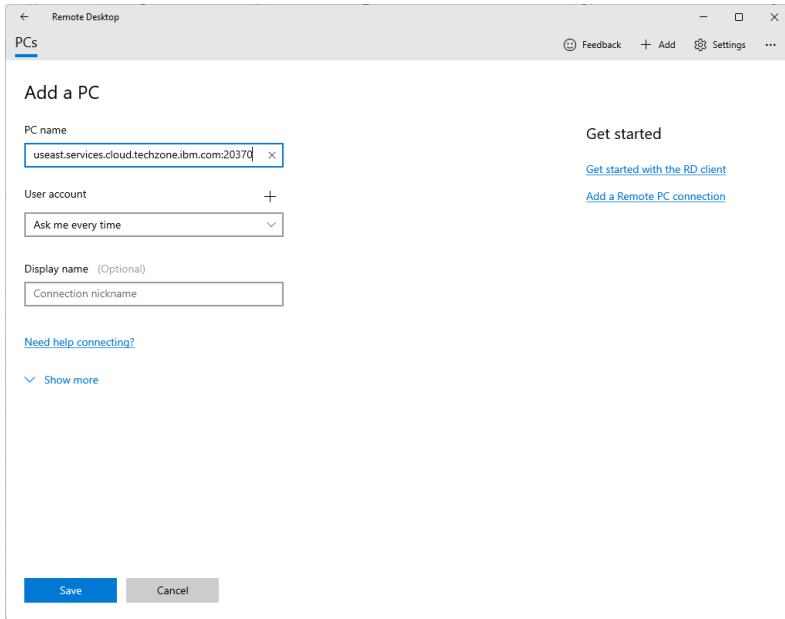
Published services

Bastion Host Remote Desktop (copy value to RDP): <useast.services.cloud.techzone.ibm.com:20370>

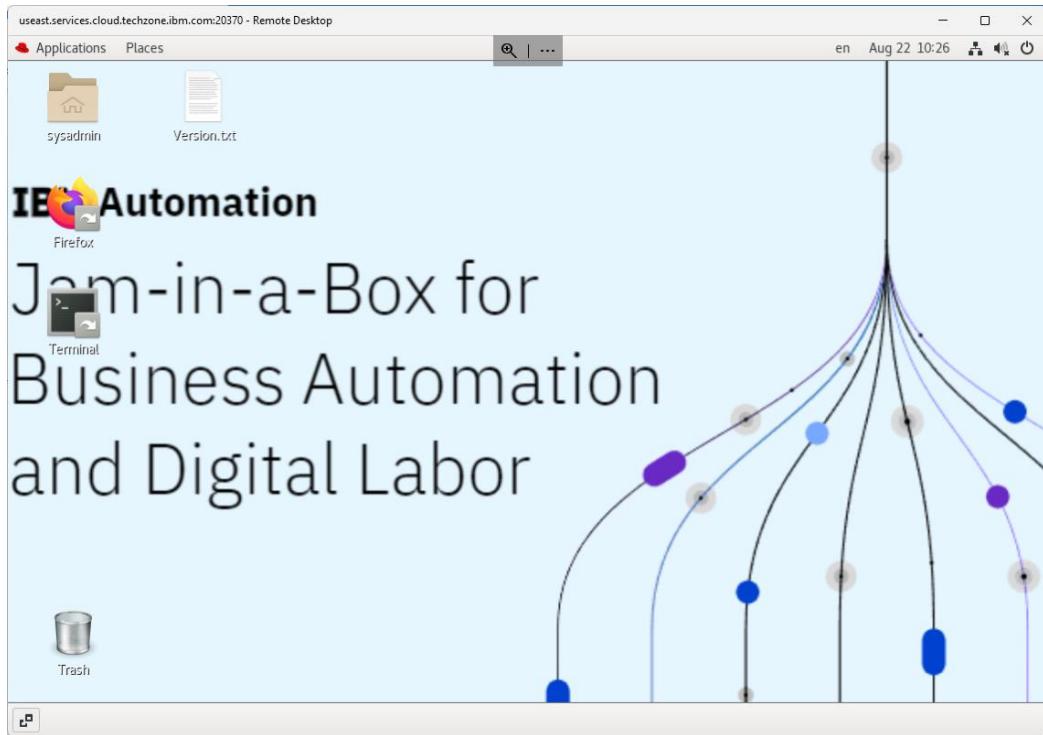
Purpose

Purpose	Opportunity ID(s)
Education	
Opportunity Product(s)	Opportunity description
Customer(s)	To learn ADP

- _3. Open the remote desktop app and paste the rdp url to the hostname as shown below

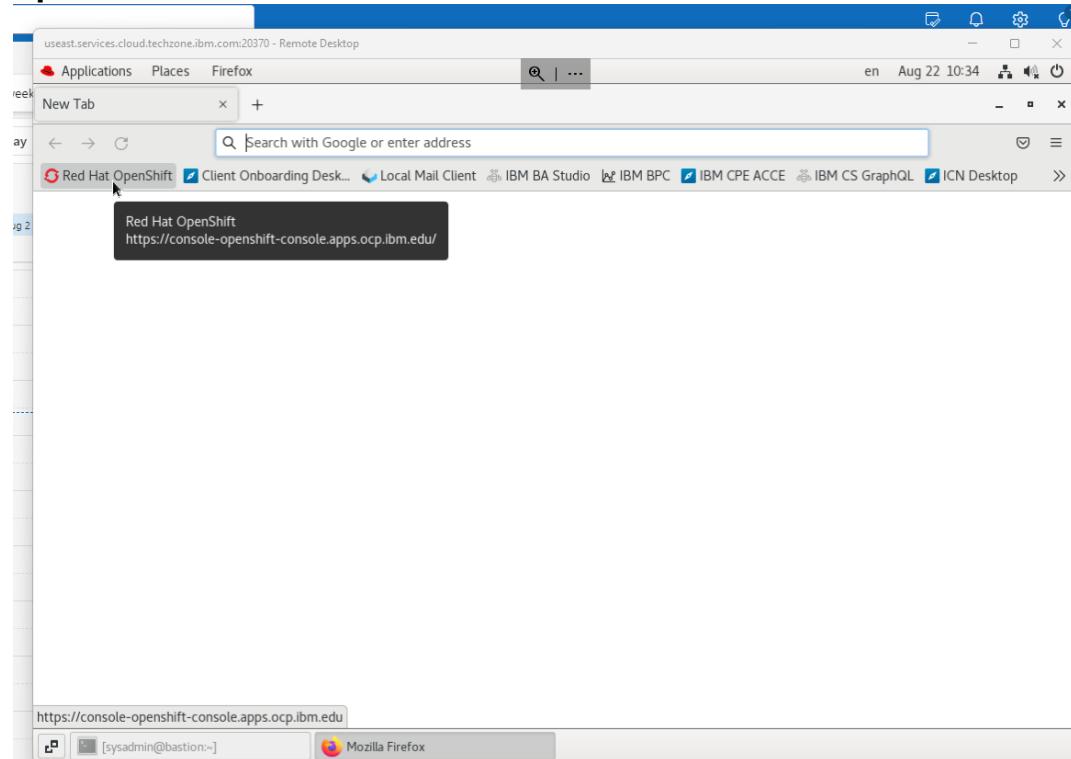


- _4. Logon using the user and password that you found on the tile when you reserved an environment. This will take you to the bastion host screen.



2.1 Preparing your Jam-in-a-Box environment

- _5. Open the **Firefox** browser from the desktop and click the first bookmark **Red Hat OpenShift Console.**



- _6. Use the populated screen for the credentials and continue to logon.

A screenshot of the Red Hat OpenShift login interface. The header features the Red Hat OpenShift logo. The main area is titled "Log in to your account". It contains two input fields: "Username" with the value "ocadmin" and "Password" with masked input. A blue "Log in" button is at the bottom. At the very bottom of the page, there is a footer message "Welcome to Red Hat OpenShift".

_7. Once logged in, on the top right corner of the window, click the **dropdown** for the **ocadmin** and select the **Copy login command**.

The screenshot shows the Red Hat OpenShift 'Overview' page. At the top right, there is a user dropdown menu for 'ocadmin'. A context menu is open over the 'Copy login command' option, listing 'Copy login command', 'User Preferences', and 'Log out'. The main interface shows cluster details like API address, ID, and version, along with status indicators for Cluster, Control Plane, Operators, and Insights. The Insights section shows it is disabled. The Activity section lists recent events such as cluster updates and server status changes.

_8. After potentially logging in again and then clicking on **Display Token** on the page that opened, copy the command below **Log in with this token** (line highlighted below).

The screenshot shows the 'Display Token' page from 'oauth Openshift.apps.ocp.ibm'. It displays an API token and a command-line login command. The command is highlighted with a red box.

```
oc login --token=sha256~aTxjUbzni2oYr1kj79ND926q7R_7hMFQ9zQjis00M0k --server=https://api.ocp.ibm.edu:6443
```

Your API token is
sha256~aTxjUbzni2oYr1kj79ND926q7R_7hMFQ9zQjis00M0k

Log in with this token

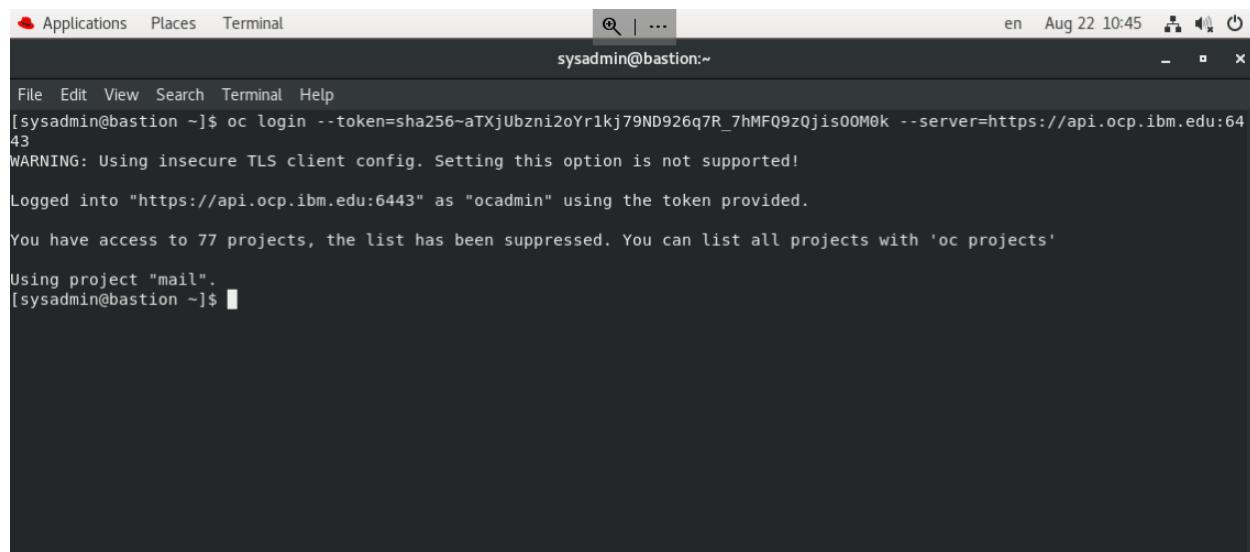
```
oc login --token=sha256~aTxjUbzni2oYr1kj79ND926q7R_7hMFQ9zQjis00M0k --server=https://api.ocp.ibm.edu:6443
```

Use this token directly against the API

```
curl -H "Authorization: Bearer sha256~aTxjUbzni2oYr1kj79ND926q7R_7hMFQ9zQjis00M0k" "https://api.ocp.ibm.edu:6443/apis/user.openshift.io/v1/users/~"
```

[Request another token](#)

_9. Open a terminal window from the desktop and paste the copied line and hit enter.

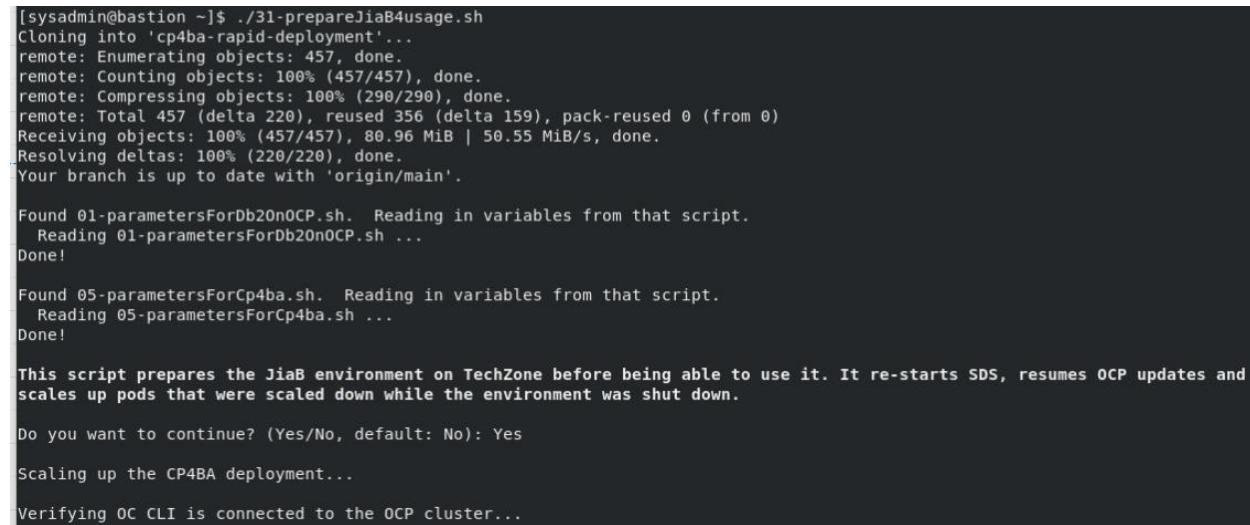


```
[sysadmin@bastion ~]$ oc login --token=sha256~aTXjUbzni2oYr1kj79ND926q7R_7hMFQ9zQjis00M0k --server=https://api.ocp.ibm.edu:6443
WARNING: Using insecure TLS client config. Setting this option is not supported!
Logged into "https://api.ocp.ibm.edu:6443" as "ocadmin" using the token provided.

You have access to 77 projects, the list has been suppressed. You can list all projects with 'oc projects'

Using project "mail".
[sysadmin@bastion ~]$
```

_10. Now type **./31-prepareJiaB4usage.sh** (including the dot at the beginning) and press **Enter** to execute the script that prepares the environment for usage.



```
[sysadmin@bastion ~]$ ./31-prepareJiaB4usage.sh
Cloning into 'cp4ba-rapid-deployment'...
remote: Enumerating objects: 457, done.
remote: Counting objects: 100% (457/457), done.
remote: Compressing objects: 100% (290/290), done.
remote: Total 457 (delta 220), reused 356 (delta 159), pack-reused 0 (from 0)
Receiving objects: 100% (457/457), 80.96 MiB | 50.55 MiB/s, done.
Resolving deltas: 100% (220/220), done.
Your branch is up to date with 'origin/main'.

Found 01-parametersForDb20nOCP.sh. Reading in variables from that script.
  Reading 01-parametersForDb20nOCP.sh ...
Done!

Found 05-parametersForCp4ba.sh. Reading in variables from that script.
  Reading 05-parametersForCp4ba.sh ...
Done!

This script prepares the JiaB environment on TechZone before being able to use it. It re-starts SDS, resumes OCP updates and scales up pods that were scaled down while the environment was shut down.

Do you want to continue? (Yes/No, default: No): Yes
Scaling up the CP4BA deployment...

Verifying OC CLI is connected to the OCP cluster...
```

When the script asks you if you want to continue, enter **y** (or Y or Yes or YES) and hit enter. Allow the script to complete and continue.

_11. Click on the fourth bookmark **IBM BA Studio** in Firefox. This would take you to the Cloud Pak login page.



It might happen that:

- *the login does not work or shows "Error 502 - Bad Gateway". In this case, please wait for some more time (about 15 minutes), then the log-in should work and the requested page is shown.*
- *the login results in "404 Page not found" error. In this case, please wait for some more time, then the log-in should work and the requested page is shown.*

These issues are the result of restarting some pods, which may take a different amount of time depending on the resources available on TechZone.

_13. On the login page, if not already selected, choose **Enterprise LDAP** under **Log in with**. Select the entry for **cp4badmin** from the list of saved logins that shows up as soon as you click into the Username field. This will also populate the password. Finally, **click Log in**.

_14. You will be presented with the “Welcome! Let’s get started” screen. In case you are offered to take a tour, **click the 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

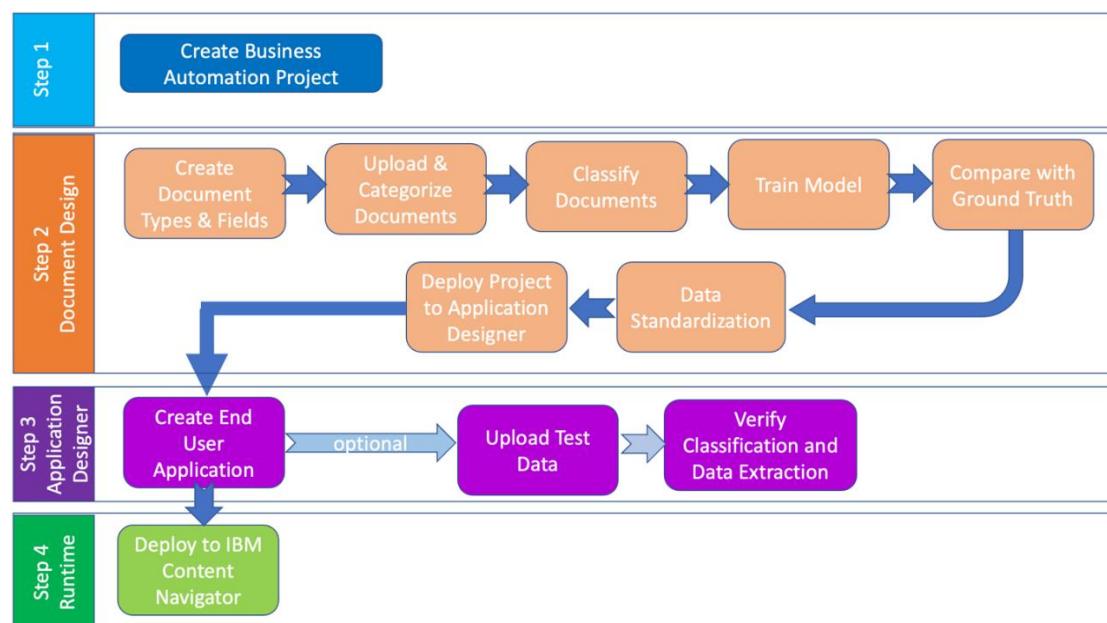
The application that you build uses AI and deep learning to automatically detect, extract, and standardize the data in all your documents. Any anomalies are flagged according to your customized model and the priority that you set so that your document processing user can correct issues before the documents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



Step1 – 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 ADP 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 Lab for Business Automation Application.

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

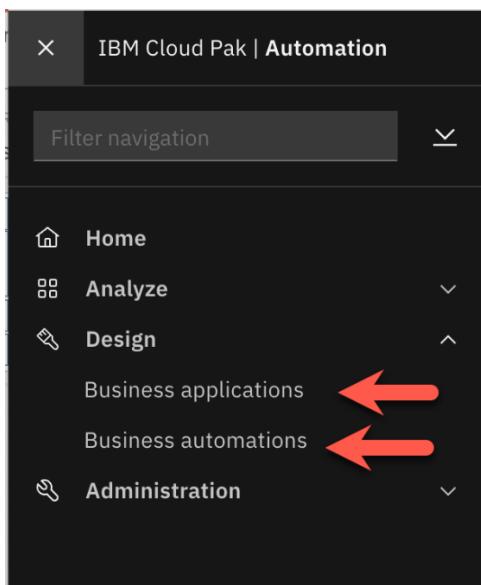
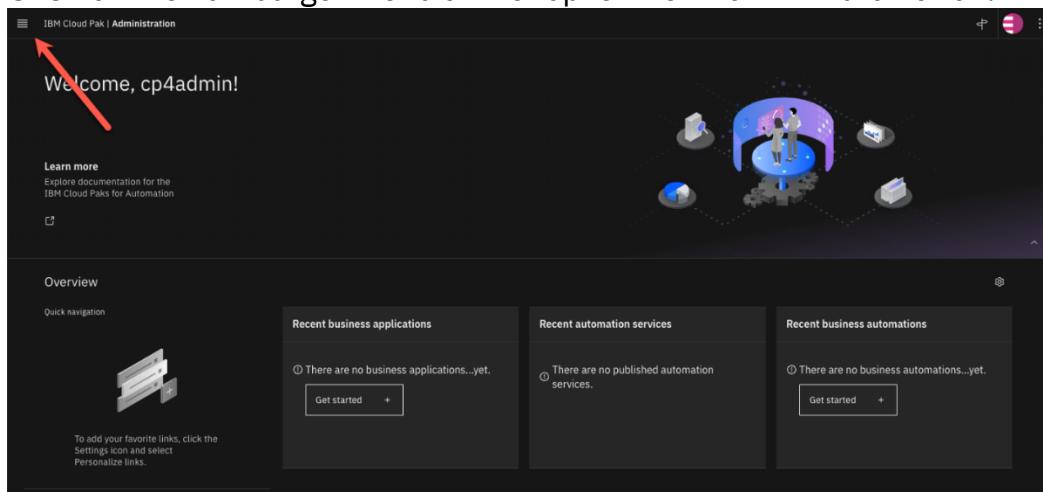
Step 1

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 access to the designer of the Document Processing configuration of the document classes, and **Business applications** provides access to the designer for 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 be called and reused 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.

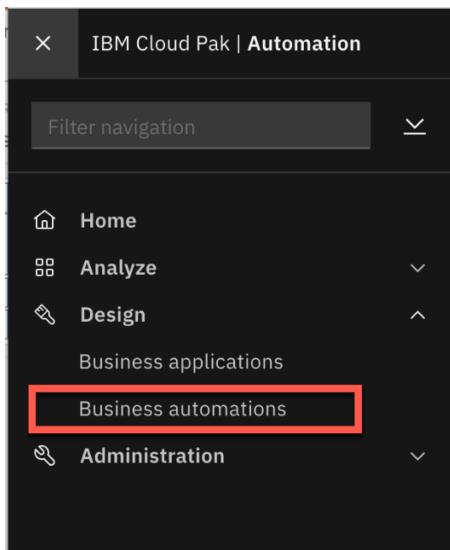
The Document Processing Designer combines an intuitive interface with a set of AI and deep learning tools that identify and learn the document types that matter to an organization. For each document type, you designate which pieces of information to extract as data for that document to be used by downstream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

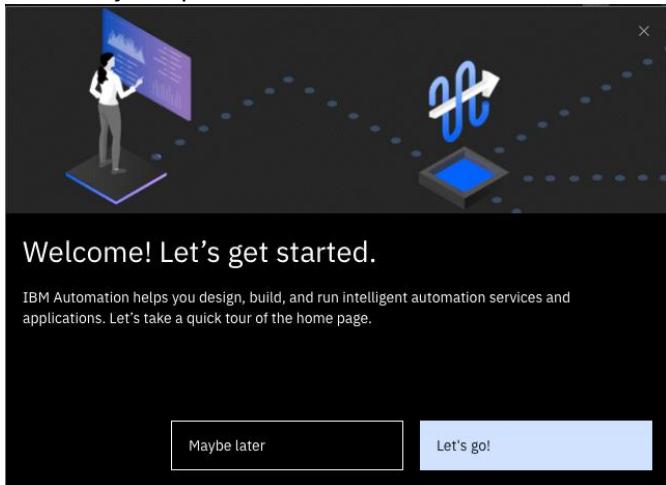
We will start with the Business Automations. Once logged in to the IBM Automation Server, you should see the Welcome screen.



2. Click on **Drop down arrow** next to Design then **Select Business automations**.



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



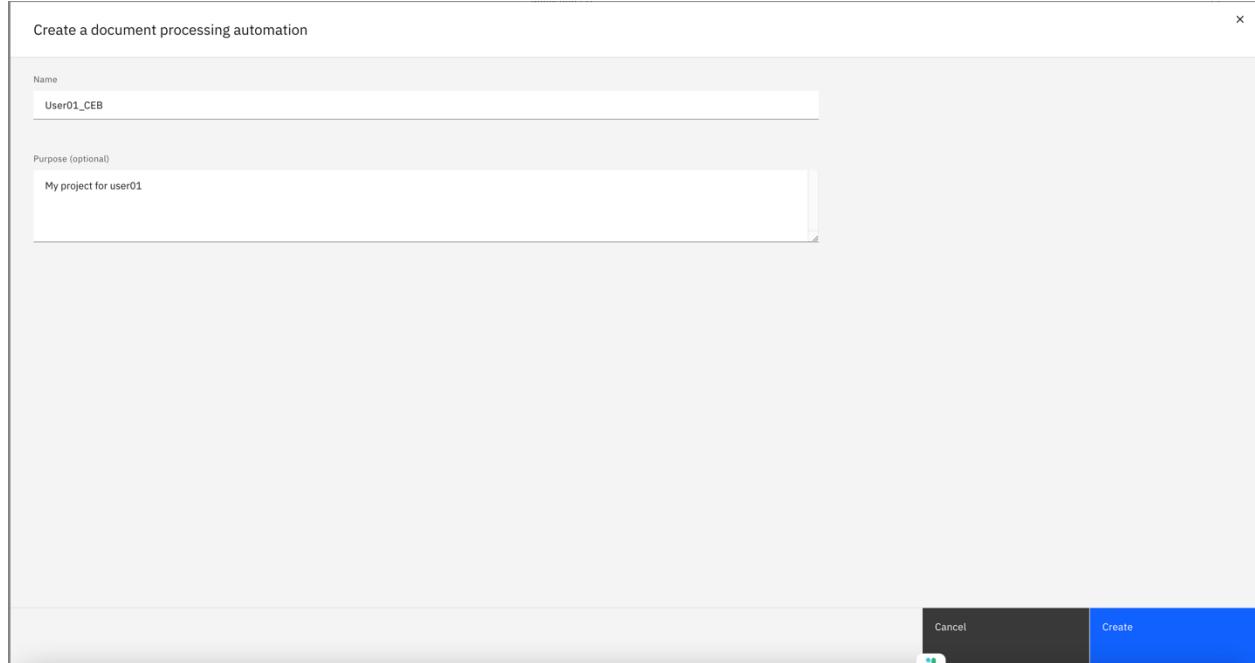
Then following screen appears.

The screenshot shows the 'Business automations' page in the IBM Cloud Pak | Automation interface. At the top, there is a navigation bar with a menu icon and 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, there are two main buttons: 'Create' (highlighted with a blue background) and 'Import'. Underneath these buttons is a section titled 'Published automation services' with a right-pointing arrow. Further down, there is a list of categories: '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**.

The screenshot shows the 'Create' dropdown menu from the previous interface. A red arrow points to the 'Create' button, which is highlighted with a blue background. The dropdown menu lists several options: 'Decision automations', 'Document processing automations' (which is highlighted with a red border), 'Workflow', and 'External'. Below the dropdown, the main content area shows the same 'Business automations' page as the first screenshot, with the 'Document processing' category selected and its details visible.

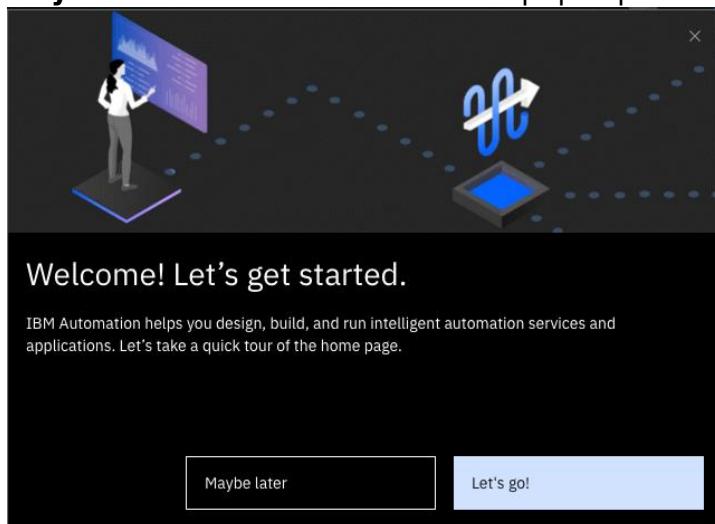
- _8. In the Create a document processing automation window **enter a name** for the project. Optionally, enter a purpose.



- _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 interface for a project named "User01_CEB". The top navigation bar includes "IBM Cloud Pak", "Business automations / User01_CEB", and three icons. Below the navigation is a horizontal menu with "Build" (underlined), "Enrich", and "Configure". The main content area is divided into five sections: "Document types and samples", "Classification model", "Extraction model", "Data standardization", and "Document retention". Each section displays status (Ready or Not ready), type counts, and accuracy percentages. In the top right corner, a yellow "Service warning" box is displayed with the message: "To resolve this, you must enter valid credentials in the Git server configuration dialog, under the Configuration tab." It also includes a "Share" button and a "Last shared" timestamp.

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

On the top right, you initially see a yellow Service warning. This is because in your environment ADP is not yet connected to a Git repository. Close this warning, you will take care of it in section 4.1.3.

Once closed, the **Share** and **Version / Deploy** buttons will be completely visible.



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: Here you will teach the system how to recognize the different document types.

Extraction model: Here you 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

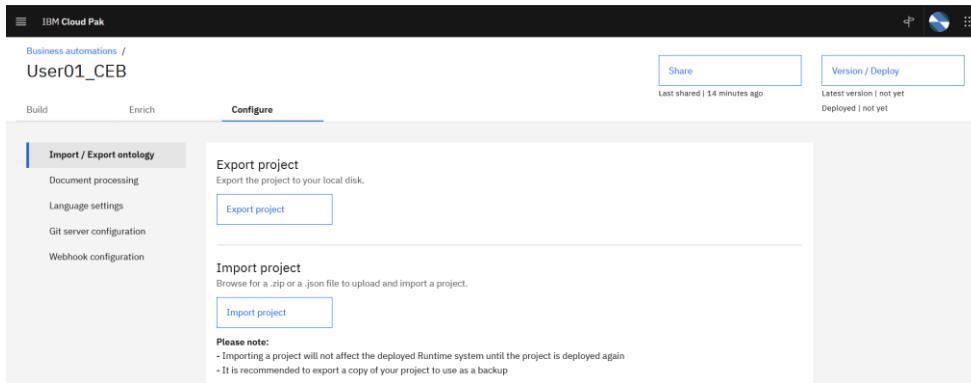
- _3. Click on <**your project name**> in the bread crumb trail at the top to go back to the Enrich tab.

4.1.3 Configure Tab

- _4. Click on **Configure** tab

This is where we can configure other operational aspects of the project.

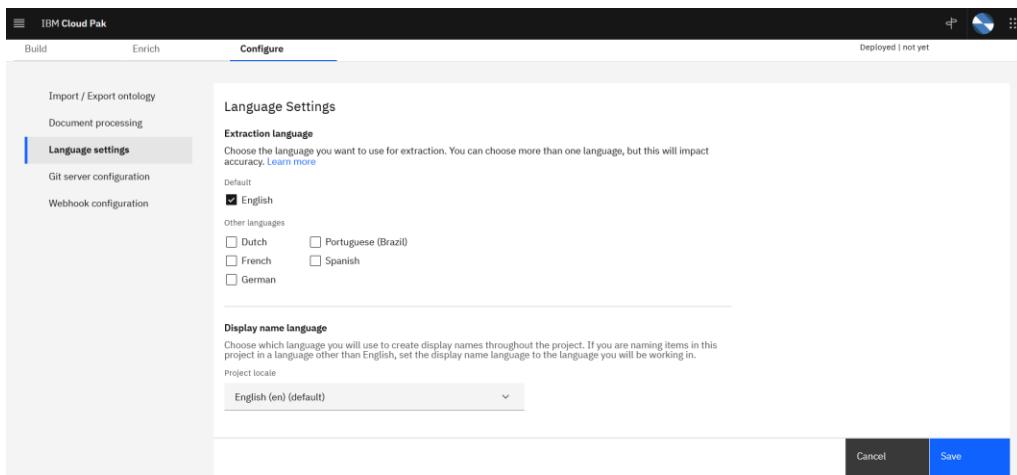
On the default tab **Import / Export project**, the **Export project** creates a .zip file that contains the document types, field types and enrichments, which you can use to start training with new sample files. You can also decide to include the training model and the sample training files in your export if you want to move your entire project to a new instance of Document Processing for example. You can import a project by clicking **Import project** selecting 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.



On the **Language settings** tab under **Extraction language**, you 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.



On the **Git server configuration** tab, 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.

Fill out the form with the following values:

- **Git vendor:** Gitea
- **Git server organization URL:** `https://simple-gitea-gitea.apps.ocp.ibm.edu/Automation-SWAT`
- **Git server REST API URL:** `https://simple-gitea-gitea.apps.ocp.ibm.edu/api/v1`
- **Username:** select **cp4badmin** from the list of saved logins
- **Type of credentials:** Password
- **Credentials:** will be auto populated when you select cp4badmin from the list of saved logins

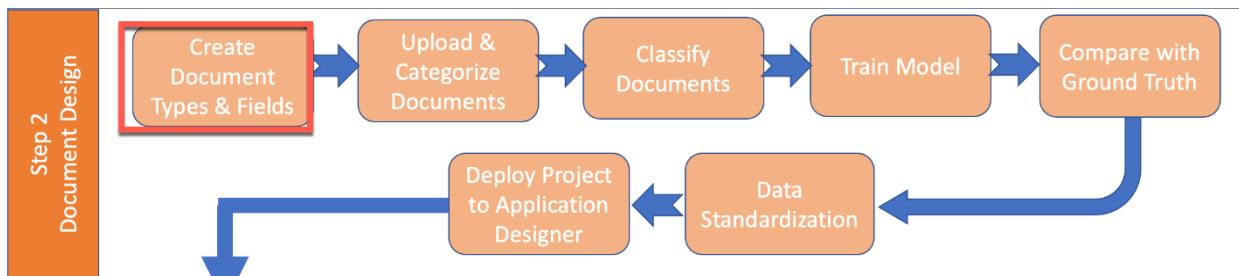
First **click** on the “**Test**” button, which should result in a **Test connection successful** message being shown in green.

Once successful, click on the “**Save**” button, this should also succeed.

After that in the top right corner **click** on the “**Share**” button. This is required to be able to create a version later.

The screenshot shows the 'Configure' tab selected in the top navigation bar. On the left, a sidebar lists 'Import / Export project', 'Document processing', 'Language settings', 'Git server configuration' (which is highlighted), and 'Webhook configuration'. The main content area displays a message: 'In order to share, version and deploy, you need to establish a connection to your organization's Git server.' Below this, there are fields for 'Git vendor' (set to 'Gitea'), 'Git server organization URL' (set to '`https://simple-gitea-gitea.apps.ocp.ibm.edu/Automation-SWAT`'), 'Git server REST API URL' (set to '`https://simple-gitea-gitea.apps.ocp.ibm.edu/api/v1`'), 'Username' (set to 'cp4badmin'), 'Type of credentials' (radio button selected for 'Password'), and 'Credentials' (represented by a redacted password field). At the bottom are two buttons: 'Test' (in a light blue box) and 'Save' (in a dark blue box).

5 Configure a Wage and Tax document type

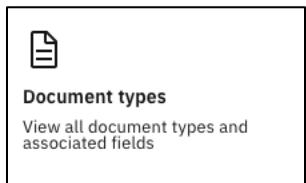


Before we use the guided configuration, you will configure some additional document types and fields used to extract data prior to uploading sample documents.

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 the **Enrich** tab
- _2. Click on **Document types**



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

- _3. Click on the **Create document type +** button in the top right corner



- _4. 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 and remove the spaces. There’s no need to add description in this lab unless you want to.

Add document type X

Display name 12/50

This is the name that will show up for you in the system. You can use characters from any language.

Symbolic name 10/50

This name will be used to identify the document type in the code.

Classification confidence threshold %
 - +

Set a confidence level to be aware of documents that fall under the desired threshold. Documents under this threshold will show a warning.

Description (optional) 0/512

Fixed format ⓘ

Feedback documents ⓘ

Percentage of corrected documents to use in retraining ⓘ
 - +

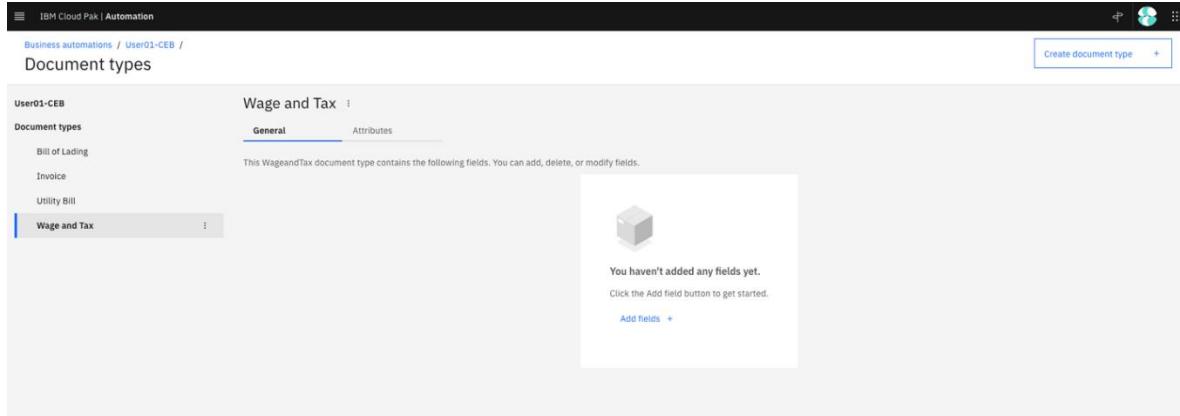
Cancel Add



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.

5. Click the Add button

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

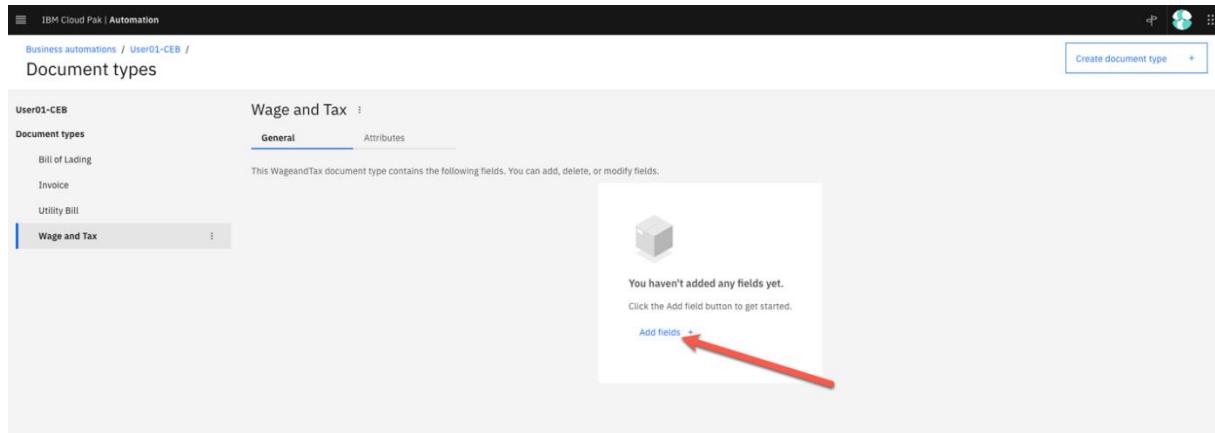


_6. Select your **Wage and Tax** doc type. On the right, you should see an empty table of fields.

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 +



Enter the following values under the **General Settings** header

The screenshot shows the 'Create field' interface in IBM Cloud Pak | Automation. The 'General' tab is selected. The 'Display name' field contains 'Ex. Employee's name, Le nom de l'employé'. A red error message 'This is a required field' is displayed below it. The 'Field type' dropdown is set to 'sys:String'. The 'Aliases' section contains 'Enter an alternative name'.

- Display name: **Federal Income Tax Withheld**
- Field type:
 - **Sys:Decimal**
- This field is required: **Yes**
- In Aliases enter other possible names. Case and punctuation are very important 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: In the second case, the number two has a period after it!

You should now see the following:

The screenshot shows the 'Create field' interface in IBM Cloud Pak | Automation. The 'General' tab is selected. The 'Display name' field contains 'Federal Income Tax Withheld'. A note below it says 'This is the name that will show up for you in the system. You can use characters from any language.' The 'Field type' dropdown is set to 'sys:Decimal'. The 'Aliases' section contains 'Enter an alternative name' and two entries: '2 Federal income tax withheld' and '2. Federal income tax'.

_2. 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.

_3. 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.

_4. 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 in the **Value format** section.

The screenshot shows the 'Value settings' tab of a document type configuration page. The 'Converters' section is highlighted with a red box, and the 'Edit' button next to it is also highlighted.

_5. Click on Converters tab then click on the blue Add converter + button

The screenshot shows the 'Text value format' screen with the 'Converters' tab selected. A red box highlights the 'Converters(0)' tab, and a blue 'Add converter +' button is visible below it.

_6. 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

- _7. You will then be presented with the 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. Refer [Enrichments-Converters](#) for more details.



Note: For Decimal cleans values such as currency to remove extra non-numeric characters and convert to the decimal format that you want. Available for the Decimal field type.

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(1)

Converter details	Test all converters
Converter name: Decimal Converter	Sample value: Enter the sample value to test the converters with
Description: Decimal Converter	Test
Type: Decimal Converter	Converted result: No results yet Converted result will be displayed here after clicking the Test button.
Decimal point: .	
Max digits after decimal point: 2	
Inherited from: sys.Decimal	

CANCEL Done

- _8. **Click Create** in the top right. Once it is created you will be taken back to the Document type page. Your screen should look like this with your first field created.

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**
- This field is 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**
- This field is 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

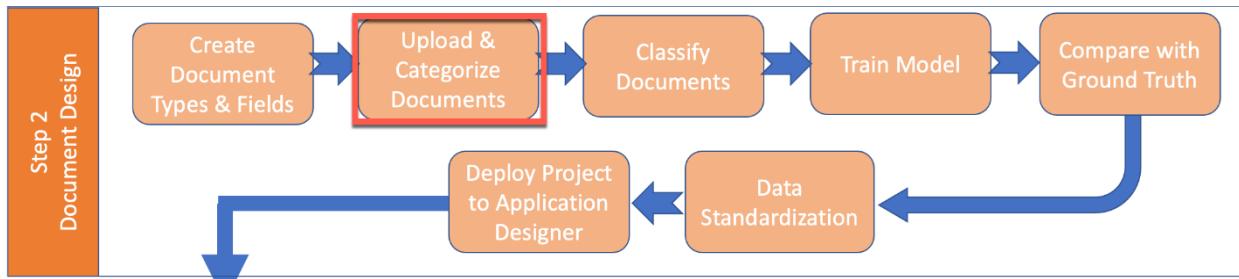
<https://www.ibm.com/docs/en/cloud-paks/cp-biz-automation/24.0.0?topic=enrichments-field-types-document-processing>

- _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 <User01_CEB>. This will take you back to the **Enrich** tab, then click on the **Build** tab.

The screenshot shows the IBM Cloud Pak interface for managing document types. The top navigation bar includes 'IBM Cloud Pak', 'Business automations / User01_CEB /', and a 'Create document type' button. On the left, a sidebar lists 'Document types' such as 'Bill of Lading', 'Invoice', 'Utility Bill', and 'Wage and Tax'. The 'Wage and Tax' item is selected and highlighted in blue. The main content area is titled 'Wage and Tax' and shows a table of fields. The table has columns for 'Name', 'Type', 'Required', and 'Sensitive'. Fields listed include 'Employee Name and Address' (String), 'Employee Social Security Number' (SocialSecurityNumber), 'Employer Identification Number' (String), 'Employers Name and Address' (String), 'Federal Income Tax Withheld' (Decimal), 'Social Security Wages' (Decimal), and 'Wages Tips Other Compensation' (Decimal). A search bar and an 'Add fields' button are also visible.

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	true	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 part of the lab, we will teach the system to recognize the various document types on your system.

In the first part of the classification section, 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 section 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 screenshot below.

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 is titled 'Document types and samples'. It contains a sub-section 'Upload sample documents to define the types of documents you want the system to process.' Below this, there are five categories: 'Classification model', 'Extraction model', 'Data standardization', and 'Document retention'. Each category has a status indicator (Ready, Retrain, Not ready), the number of types (4, 3, 3, 4 respectively), and accuracy percentages (100%, 96%, N/A, 100%). A red box highlights the 'Document types and samples' section at the top.

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

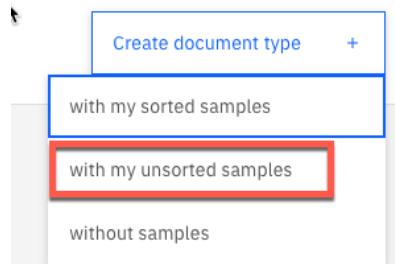
The screenshot shows a dropdown menu triggered by a blue button labeled 'Create document type +'. The menu contains three items: 'with my sorted samples', 'with my unsorted samples', and 'without samples'. The first item, 'with my sorted samples', is highlighted with a blue border.

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.				
<input type="text"/> Search sample documents Upload 				
<input type="checkbox"/>	Document name	Status	Date Added	Added by
<input type="checkbox"/>	Mortgage Agreement1.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	Mortgage Agreement2.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	Mortgage Agreement3.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	Mortgage Agreement4.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	Mortgage Agreement5.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	TR_FW2_1001_0000_PS.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	TR_FW2_2000_0000_PS.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	TR_FW2_3000_0000_PS.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	TR_FW2_3001_0000_PS.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin
<input type="checkbox"/>	TR_FW2_4000_0000_PS.pdf	✓ Ready	Feb 22, 2024 11:53 AM	cp4admin

Items per page 10

1 - 10 of 12 items

1 of 2 pages

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

The screenshot shows a web-based interface for document classification. At the top, there's a navigation bar with 'IBM Cloud Pak' and a breadcrumb trail: 'Business automation / User01_CIB / Document types and samples /'. Below the navigation is a section titled 'Create document types' with two options: 'Upload unsorted documents' (selected) and 'Review categories'. A note below says '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.' On the left is a search bar labeled 'Search sample documents'. The main area is a table with columns: 'Document name', 'Status', 'Date Added', and 'Added by'. The table lists several PDF files: 'Mortgage Agreement1.pdf', 'Mortgage Agreement2.pdf', 'Mortgage Agreement3.pdf', 'Mortgage Agreement4.pdf', 'Mortgage Agreement5.pdf', 'TR_FW2_1001_0000_P5.pdf', 'TR_FW2_2000_0000_P5.pdf', 'TR_FW2_3000_0000_P5.pdf', 'TR_FW2_3001_0000_P5.pdf', and 'TR_FW2_4000_0000_P5.pdf'. All documents are marked as 'Ready' and were added on Feb 22, 2024, at 11:53 AM by 'cp4admin'. At the bottom of the table are buttons for 'Items per page' (set to 10), a page number indicator ('1 - 10 of 12 items'), and a 'Categorize' button.



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.

The order of the categories shown in the screenshots below may differ from the order in your environment.

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. On the left, there's a sidebar with 'Categories (3)' (Category 1 is selected), 'Document types (4)' (Bill of Lading, Invoice, Utility Bill, Wage and Tax), and a 'Review categories' button. The main area shows 'Category 1 sample documents (2)'. It has a search bar and a table with columns: Document name, Status, Date Added, and Added by. Two PDF files are listed: 'UBILLCable_081_1_1.pdf' and 'UBILLCable_082_1_1.pdf', both marked as 'Ready' and added on Feb 22, 2024, by cp4admin.

This screenshot shows the 'Create document types' interface in IBM Cloud Pak. The sidebar shows 'Categories (3)' (Category 2 is selected), 'Document types (4)' (Bill of Lading, Invoice, Utility Bill, Wage and Tax), and a 'Review categories' button. The main area shows 'Category 2 sample documents (5)'. It has a search bar and a table with columns: Document name, Status, Date Added, and Added by. Five PDF files are listed: 'Mortgage Agreement1.pdf', 'Mortgage Agreement2.pdf', 'Mortgage Agreement3.pdf', 'Mortgage Agreement4.pdf', and 'Mortgage Agreement5.pdf', all marked as 'Ready' and added on Feb 22, 2024, by cp4admin.



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.

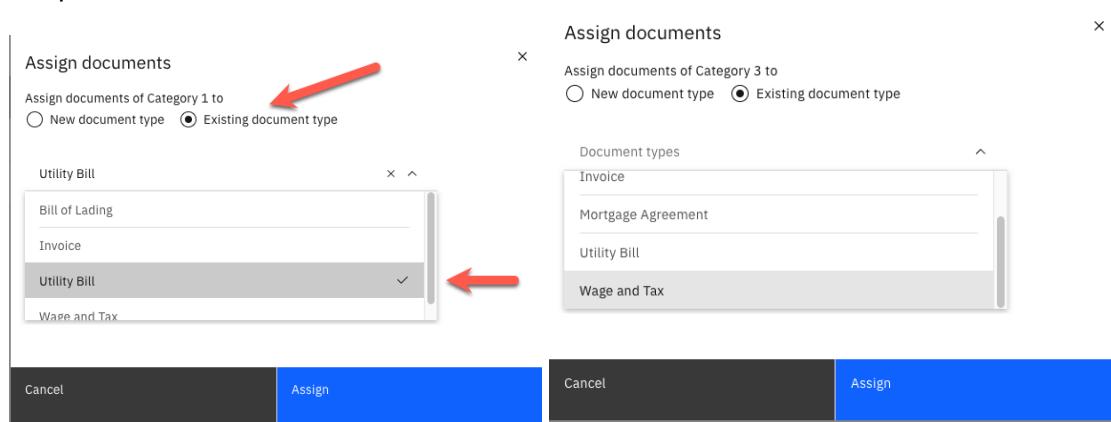
For each of the three categories perform the following steps:

- _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. If the documents are either of type **Utility Bill** or **Wage and Tax**:

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



Click Assign to close the dialog box.

If the documents are of type **Mortgage Agreement**:

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

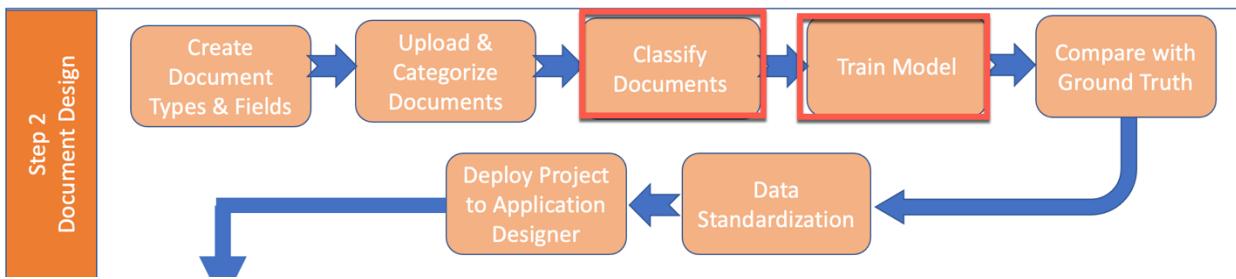
Enter Mortgage Agreement in the field

The dialog box is titled 'Assign documents' and has the sub-instruction 'Assign documents of Category 2 to'. It contains two radio buttons: 'New document type' (checked) and 'Existing document type' (unchecked). Below this is a section for 'Document type display name' with a character limit of 18/50, containing the text 'Mortgage Agreement'. Below this is a note: 'This is the name that will show up for you in the system. You can use characters from any language.' Below this is a section for 'Document type symbolic name' with a character limit of 17/50, containing the text 'MortgageAgreement'. At the bottom are 'Cancel' and 'Assign' buttons.

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

_ 10. Click the **Finish** button in the top right corner

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 document 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 get recognized 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 bread crumb trail to return to the start page

_2. Click anywhere in the **Classification model** line

The screenshot shows the IBM Cloud Pak Business Automation interface. At the top, there's a navigation bar with 'IBM Cloud Pak' and a user icon. Below it, the breadcrumb trail shows 'Business automation / User01_CEB'. On the right, there are 'Share' and 'Version / Deploy' buttons with status information: 'Last shared 1 hour ago' and 'Latest version not yet Deployed'.

The main content area has tabs: 'Build' (selected), 'Enrich', and 'Configure'. Under 'Build', there are several sections:

- Document types and samples**: Shows 5 types and 20 samples on average.
- Classification model** (highlighted with a red box): Shows 3 types trained with 100% accuracy. It includes a sub-instruction: 'Train the model to classify your documents.'
- Extraction model**: Shows 3 types trained with 96% accuracy. It includes a sub-instruction: 'Train the model to extract the data from your documents.'
- Data standardization**: Shows 'Not ready'.
- Document retention**: Shows 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 document types 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
Last trained: a day ago
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. Include at least 1 document in the test set to view training results.

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

Classification model
Last trained: a day ago
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 60/40 split to the test set. Click **Auto generate split**

Document types

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

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](#)

Training/test ratio in %
100/0

Test set (0) 0% of total samples

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 60/40 – that is, 60% of the available sample documents should be used for training, and we will validate the training results with 40% of the sample documents. This split is only a suggestion, and we can adjust it, but 60/40 is a good starting point.

Classification model

Accuracy 84.8%

Training/test ratio in %
60/40

Test set (2) 40% of total samples

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_005_2_1.pdf	Bill of Lading	Correct	96.06%
BOL_009_2_1.1.pdf	Bill of Lading	Correct	91.63%
BOL_015_2_1.1.pdf	Bill of Lading	Correct	95.23%
BOL_027_2_1.1.pdf	Bill of Lading	Correct	95.68%
BOL_041_2_1.1.pdf	Bill of Lading	Correct	94.72%
BOL_051_2_1.1.pdf	Bill of Lading	Correct	93.37%
BOL_054_2_1.1.pdf	Bill of Lading	Correct	96.14%

7. Close the green notification. Click on each of the document types. Notice the confidence levels. You can notice either or both Mortgage Agreement or 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** page. 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**. Wait until the status for all documents is Ready.
- _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' page for the 'Wage and Tax' category. On the left, there's a sidebar with 'Document types' listed: Bill of Lading (28 samples), Invoice (26 samples), Mortgage Agreement (1 sample), Utility Bill (23 samples), and Wage and Tax (5 samples). The 'Wage and Tax' category is selected. The main area displays a table titled 'Wage and Tax sample documents (5)'. The table has columns: Document name, Status, Date Added, and Added by. All five documents listed are 'Ready' and were added on Sep 19, 2023, at 11:18 PM, by user cp4admin. A red box highlights the 'Upload' button in the top right corner of the table header.

Document name	Status	Date Added	Added by
TR_FW2_1001_0000_PS.pdf	✓ Ready	Sep 19, 2023 11:18 PM	cp4admin
TR_FW2_2000_0000_PS.pdf	✓ Ready	Sep 19, 2023 11:18 PM	cp4admin
TR_FW2_3000_0000_PS.pdf	✓ Ready	Sep 19, 2023 11:18 PM	cp4admin
TR_FW2_3001_0000_PS.pdf	✓ Ready	Sep 19, 2023 11:18 PM	cp4admin
TR_FW2_4000_0000_PS.pdf	✓ Ready	Sep 19, 2023 11:18 PM	cp4admin

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

The screenshot shows the 'Classification model' page in the IBM Cloud Pak Administration interface. At the top, it displays 'Business automations / Clandis Baker Project / Classification model' with an accuracy of 100% from an hour ago. Below this are tabs for 'Confirm inputs', 'Review samples' (which is active), 'Review training results', and 'Test trained model (optional)'. A message box indicates changes were made since the last train. The 'Document types' section lists 'Bill of Lading', 'Invoice', 'Mortgage Agreement', 'Utility Bill', and 'Wage and Tax' (selected). The 'Training set (7)' contains 70% of total samples with files: TR_FW2_1000_0001F.pdf, TR_FW2_1000_0002F.pdf, TR_FW2_2000_0000_PS.pdf, TR_FW2_2000_0001F.pdf, TR_FW2_2000_0002F.pdf, and TR_FW2_3001_0000_PS.pdf. The 'Test set (3)' contains 30% of total samples with files: TR_FW2_1001_0000_PS.pdf, TR_FW2_3000_0000_PS.pdf, and TR_FW2_4000_0001F.pdf. A red box highlights the 'Auto generate 70/30 split' button at the top right of the training/test set area.

_9. Click Train button and wait until the training is complete

_10. Now look at the confidence score for **Wage and Tax**. They should have improved considerably compared to before you added new documents.

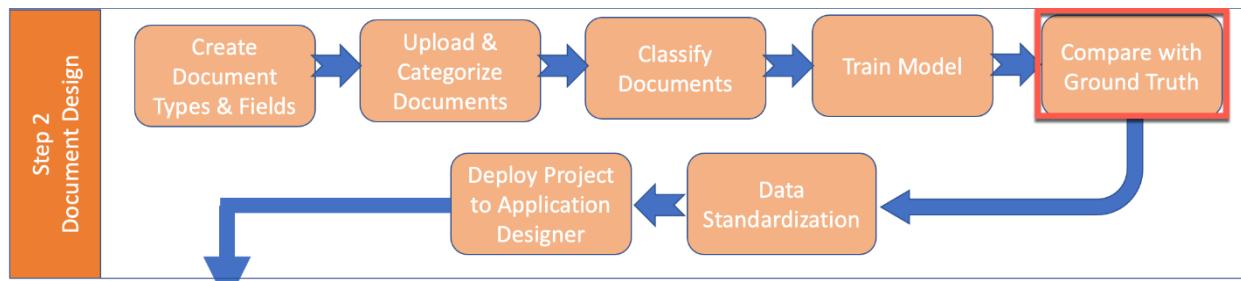
_11. Click Next and then Click Done

7.1.2 Option 2 – Review all uploaded samples

As pointed out before, the quality of the sample documents determines the quality of the results. Therefore in general:

- 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



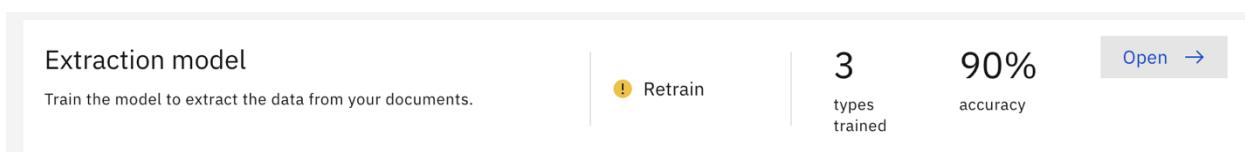
At this point, we have defined a document type, told the system which fields we want off the document and trained the system on how to recognize (classify) the document. In the Data Extraction portion of the lab, we will upload new Wage and Tax documents to Document Processing and see how our earlier 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 be reset to Retrain if ADP detects something may have changed. This is just a reminder that if you indeed changed something, you may benefit from retraining the model.



- _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.

The screenshot shows the 'Extraction model' page in the IBM Cloud Pak interface. At the top, there are tabs: 'Review samples' (selected), 'Add fields', 'Teach model', 'Review training results', and 'Test model' (Optional). Below the tabs, a message box says: 'Please make sure you have at least 1 reviewed document to train the model.' A note below it says: 'Review your training and test sets. The test set is used to generate the model training results. Learn more'. Under 'Document types', 'Wage and Tax' is selected. In the center, there are two sections: 'Training set (13)' containing 65% of total samples (documents: TR_FW2_1000_0001F.pdf, TR_FW2_1000_0002F.pdf, TR_FW2_1000_0003F.pdf, TR_FW2_1000_0004F.pdf, TR_FW2_1000_0005F.pdf, TR_FW2_1001_0000_P5.pdf); and 'Test set (7)' containing 35% of total samples (documents: TR_FW2_2000_0001F.pdf, TR_FW2_2000_0002F.pdf, TR_FW2_3001_0000_PS.pdf, TR_FW2_4000_0000_PS.pdf, TR_FW2_4000_0002F.pdf, TR_FW2_4000_0003F.pdf). There are 'Search' bars for each section and 'Next' and 'Previous' navigation buttons between them.

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

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



You will now be on the *Add fields* step. 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* step.

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 'Extraction model' interface in IBM Cloud Pak. 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. In the center, a message says 'Please make sure you have at least 1 reviewed document to train the model.' Below it is a list of 'Wage and Tax sample documents (12)'. A red box highlights the 'Teach samples' button in the top right corner of this list area. At the bottom right of the interface are 'Upload' and 'Reanalyze' buttons.



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 'Review training results' screen for the 'TR_FW2_1001_0000_PS.pdf' document. On the left, the document content is displayed, showing a W-2 form for an employee named David Smith. On the right, a table lists extracted fields with their captured values. The table includes columns for 'Field Name', 'Value Captured', and 'Type'. Fields listed include 'Employee Social Security Number', 'Employee Name and Address', 'Employer Identification Number', and 'Social Security Wages'. A checkbox at the bottom allows marking the document as ready for training.



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.

The screenshot shows the IBM Cloud Pak Administration interface with a document titled "TR_FW2_1000_0001F.pdf". The document is a 2020 W-2 Wage and Tax Statement. The right side of the screen displays a "Field Name" and "Value Captured" table with several rows of extracted data. One row for "Federal Income Tax Withheld" is highlighted with a blue border. A modal window titled "Recommended matches" lists two entries: "Federal income tax withheld 1800.00" and "2 Federal income tax withheld 1800.00". Buttons for "Edit selection", "Dismiss", and "Save selection" are visible at the bottom of the modal.

Field Name	Value Captured
Federal Income Tax Wit...	abc Text
Recommended matches ⓘ Matches are ranked in order of confidence. Choose one and save or dismiss to draw your own.	
Field label	Field value
Federal income tax withheld 1800.00	
2 Federal income tax withheld 1800.00	
Edit selection Dismiss Seeing duplicates?	
Pending aliases View all aliases (3)	
Detected alias already exists ⓘ	
Save selection	

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.

IBM Cloud Pak | Administration

Back TR_FW2_1000_0001F.pdf | Not ready

Show detected fields Keyboard shortcuts on

Sort by: Date created

Field Name Value Captured

Federal Income Tax Withheld 1800.00

Required

Recommended matches

Matches are ranked in order of confidence. Choose one and save or dismiss to draw your own.

Field label Field value

Federal income tax 1800.00 •

2 Federal income tax 1800.00 withheld

Federal income tax 1800.00 withheld

Local income tax 500.00

Edit selection Dismiss Seeing duplicates?

Pending aliases | View all aliases (3)

Detected alias already exists ⓘ

Save selection

Employee Name and A... Required

Mark this document as ready for training. ⓘ

Previous sample Next sample

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

IBM Cloud Pak | Administration

Back TR_FW2_1001_0000_PS.pdf | Not ready

Show detected fields Keyboard shortcuts on

Sort by: Date created

Field Name Value Captured

Federal Income Tax Withheld 123456789.99

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 A... Required

Employee Social Security Number Required

Employer Identification Number Required

Employers Name and Address Required

Mark this document as ready for training. ⓘ

Previous sample Next sample

Match data underlined in blue to the selected field or draw your own boxes around data in the document.

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

_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.**

The screenshot shows the IBM Cloud Pak Administration interface with the W-2 Wage and Tax Statement document loaded. A context menu is open over the 'Federal Income Tax Withheld' field, which has a green checkmark indicating it is complete. The menu options are 'Save match' (highlighted with a red box) and 'Cancel'.

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 123456789.99				
c Employer's name, address, and ZIP code Long Lengthy Name The Corporation 56334 Full Sized Avenue Unit 1234 Minneapolis, Minnesota 55411-1234	2 Federal income tax withheld 123456789.99				
d Control number 123456 A78	3 Social security wages 123456789.99				
e Employee's first name and initial Last name Michael Robert David Smithson III 56334 Full Sized Avenue Unit 1234 Minneapolis, Minnesota 55411-1234	Suff.	4 Social security tax Captured value 123456789.99			
f Employee's address and ZIP code 15 State MN	16 State wages, tips, etc. 123456789	17 State income tax 123456789.99	18 Local wages, tips, etc. 123456789.99	19 Local income tax 123456789.99	20 Locality name ABCDEFGr
W-2 Wage and Tax Statement 2020 Department of the Treasury—Internal Revenue Service Copy 1—For State, City, or Local Tax Department					

Notice a green check mark signifies this field is complete.

The screenshot shows the IBM Cloud Pak Administration interface with the W-2 Wage and Tax Statement document loaded. A context menu is open over the 'Federal Income Tax Withheld' field, which now has a red border around its input box. The menu options are 'Save match' (highlighted with a red box) and 'Cancel'.

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 123456789.99				
c Employer's name, address, and ZIP code Long Lengthy Name The Corporation 56334 Full Sized Avenue Unit 1234 Minneapolis, Minnesota 55411-1234	2 Federal income tax withheld 123456789.99				
d Control number 123456 A78	3 Social security wages 123456789.99				
e Employee's first name and initial Last name Michael Robert David Smithson III 56334 Full Sized Avenue Unit 1234 Minneapolis, Minnesota 55411-1234	Suff.	4 Social security tax 123456789.99			
f Employee's address and ZIP code 15 State MN	16 State wages, tips, etc. 123456789	17 State income tax 123456789.99	18 Local wages, tips, etc. 123456789.99	19 Local income tax 123456789.99	20 Locality name ABCDEFGr
W-2 Wage and Tax Statement 2020 Department of the Treasury—Internal Revenue Service Copy 1—For State, City, or Local Tax Department					

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 name but missed the address. 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**”.

The screenshot shows a W-2 Wage and Tax Statement for the year 2020. The document is displayed in a browser window titled "IBM Cloud Pak | Administration". A sidebar on the right contains a "Recommended matches" section with several entries. Two entries are highlighted with green boxes: "e Employee's first name and initial" and "f Employee's address and ZIP code". Below these, two more entries are shown without boxes: "4326 Aldrich Rd Minneapolis, MN 55412" and "4326 Aldrich Rd Minneapolis, MN 55412". At the bottom of the sidebar, there is a "Save selection" button.

If you got the blue lines, you would notice that only the “e Employee’s first name and initial” have blue marks. In this case the values for name and address where 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**

This screenshot shows the same W-2 form as above, but with different selection highlights. The "Employee's first name and initial" field (labeled "e") and the "Employee's address and ZIP code" field (labeled "f") are now highlighted with red boxes. The sidebar on the right displays a table of captured fields with their corresponding values. The "Employee's first name and initial" field is listed as "e Employee's first name and initial" with the value "Benjamin P. Charles". The "Employee's address and ZIP code" field is listed as "f Employee's address and ZIP code" with the value "4326 Aldrich Rd Minneapolis, MN 55412". There is also a "Save selection" button at the bottom of the sidebar.

- _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 with the following details:

- Document View:** Shows the W-2 Wage and Tax Statement form with fields like Employee's social security number (577-22-3048), Employer identification number (14-023285), and various tax withholdings.
- Analysis Panel:** On the right, it shows a breakdown of extracted fields:
 - Employer Identification...**: etc 14-023285
 - Employers Name and A...**: etc Test and Rest Inc. 563 Stoney Brook Rd, Minneapolis, MN 55411
 - Social Security Wages**: etc 1113.33
 - Wages Tips Other Com...**: etc 18000.00
- Matched Values:** A list of recommended matches for "Wages, tips, other compensation" with values 18000.00 and 18000.00. The value 18000.00 is highlighted with a blue box.
- Action Buttons:** At the bottom right are buttons for "Edit selection", "Dismiss", and "Seeing duplicates?".
- Bottom Bar:** Includes "Pending aliases", "View all aliases (5)", "Save selection", "Previous sample", and "Next sample".

- !** _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 interface with three tax forms displayed side-by-side:

- Form W2 Wage and Tax Statement (Copy 1 - For State, City or Local Tax Department):** OMB No. 1545-0003. Fields include Employee Social Security Number (98-7469372), Employee ID number (AAHE 48), Control number (AAHE 48), Employee Name & Address (David Gomez, Top Heights Markets LLC, 363 Ave And 1st St, New York, NY 10037), and various wage and tax details.
- Form W2 Wage and Tax Statement (Copy 2-To be Filed with Employer's State, City or Local Income Tax Return):** OMB No. 1545-0008. Similar structure to Form W2, showing wages and taxes for David Gomez.
- Form C - For Employee's Records:** OMB No. 1545-0008. Shows employee information (Employee Social Security Number 98-7469372, Employee ID number AAHE 48, Control number AAHE 48) and a breakdown of wages and taxes for David Gomez.

A search bar at the top right contains the query "Wages tips Other Comp..." with a result count of 1. Below the search bar, there are buttons for "Field label (optional)" (Draw icon), "Field value" (Draw icon), and "Pending aliases" (View all aliases).

8.2 Train extraction model

We will be performing the **fast training** in this lab due not having a GPU available in the environment. A GPU is only needed in 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 model button

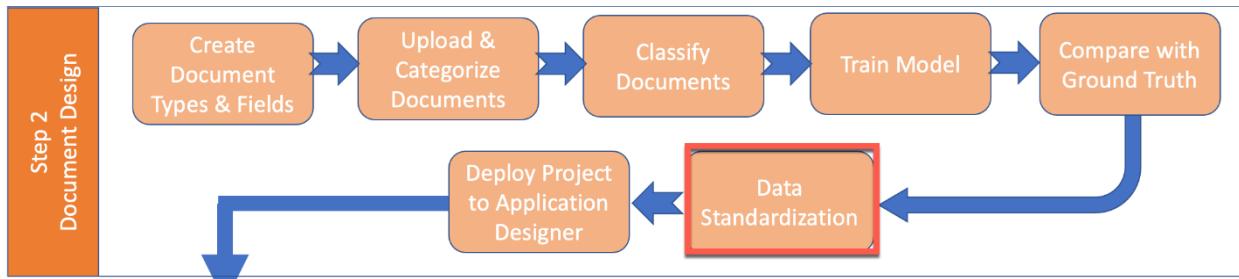
In the **Confirm training** dialog coming up, switch **Fast training** on before clicking the **Confirm** button. Then the training will take several minutes (good time for a break). If fast training is not switched on it could take days without a GPU.

The screenshot shows the "TestAdpApp / Extraction model" screen with the following interface elements:

- Top navigation: IBM Cloud Pak, Business automation / TestAdpApp / Extraction model.
- Toolbar buttons: Review samples, Add fields, Teach model, Review training results, Test model optional.
- Document types sidebar: Bill of Lading, Invoice, Mortgage Agreement, Utility Bill, Wage and Tax.
- Bill of Lading document list: BOL_001, BOL_009, BOL_015, BOL_019, BOL_026_1_1.1.pdf, BOL_027_1_1.1.pdf.
- Search bar: Document name (Search BOL_001).
- Confirm training dialog box:
 - Header: Confirm training
 - Text: Training the model can take several hours or multiple days depending on how many document types you are training. If you have other models in the queue, training starts after those models are complete.
 - Options:
 - Fast training (switched on)
 - Feedback documents (switched off)
 - Buttons: Cancel, Confirm.
- Training results table:

Date added	Ready for training
Jun 28, 2022 5:06 PM	16/16
Jun 28, 2022 5:06 PM	16/16
Jun 28, 2022 5:06 PM	16/16
Jun 28, 2022 5:07 PM	16/16

9 Data standardization

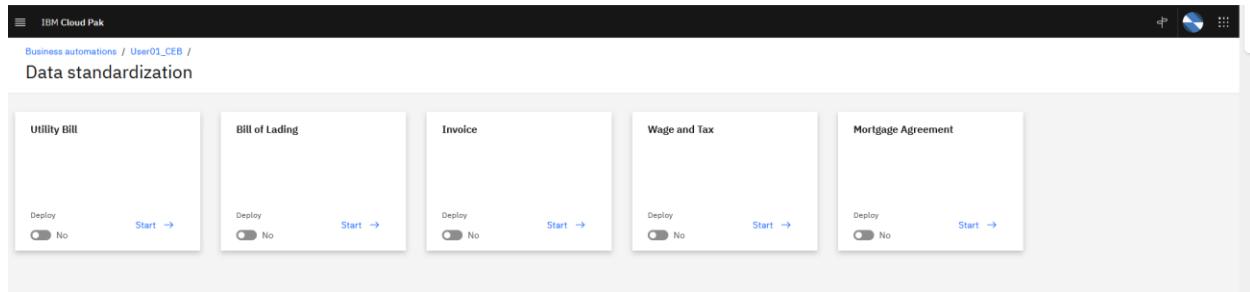


Next, we may need to standardize the data that will be presented in the user interface and how it will be stored in the FileNet repository for example. Data standardization is the process of defining attributes for a data field in a standardized way. This is done using data definitions. These definitions 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.

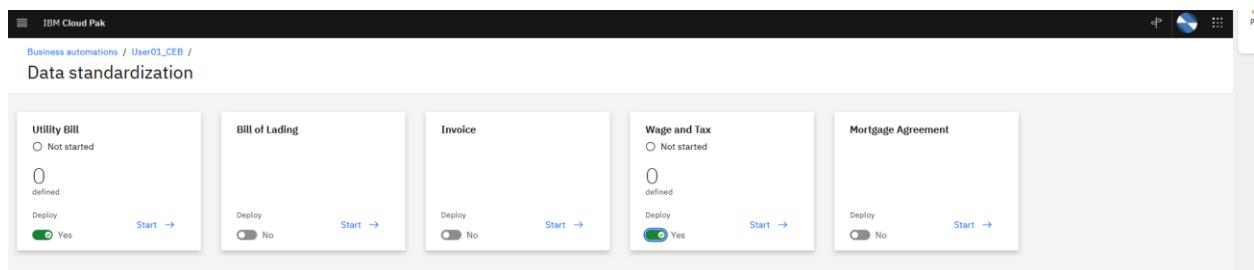
- _1. Return to the guided configuration flow and **Click** anywhere in the **Data standardization** box

Section	Status	Type	Count	Description
Document types and samples	<input checked="" type="radio"/> Ready	types	5	samples on average
Classification model	<input type="radio"/> Retrain	types trained	5	100% accuracy
Extraction model	<input checked="" type="radio"/> Ready	types trained	4	
Data standardization	<input type="radio"/> Not ready			Map fields to new or existing data definitions.
Document retention	<input checked="" type="radio"/> Ready	types reviewed	5	Determine how long you want documents to stay in your content repository.

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



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

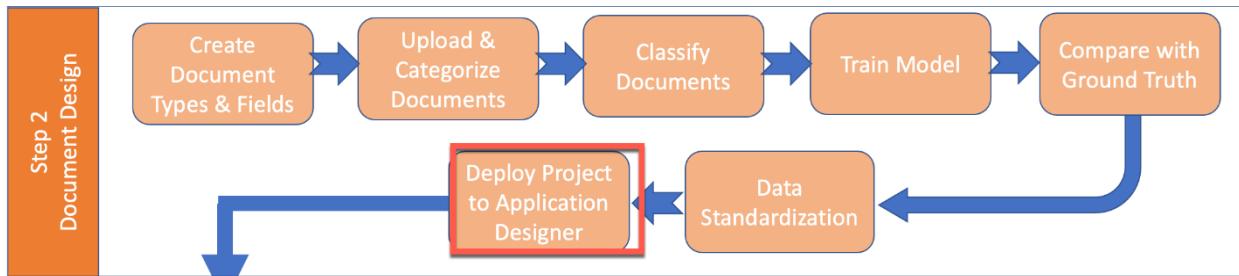


_3. Click on **Start** on either of the selected deployments

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 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.

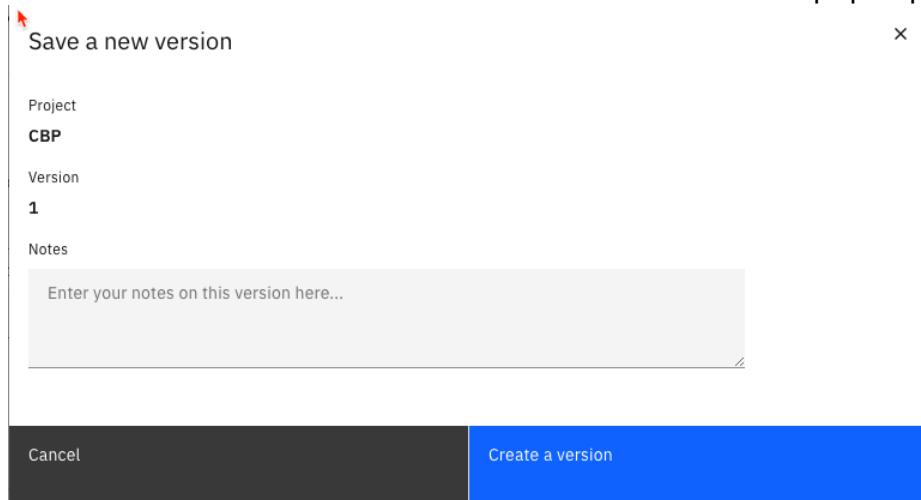
Now that we completed the configuration of the content extraction project, we need to save and deploy the design project to the application side. This will allow you to test your project using a client runtime interface.

- _1. If not already there, return to the guided home screen by clicking on your project name. Then **Click Version / Deploy**.

Document types and samples	<input checked="" type="radio"/> Ready	5 types	23 samples on average
Upload sample documents to define the types of documents you want the system to process.	→		
Classification model	<input type="radio"/> Retrain	5 types trained	100% accuracy
Train the model to classify your documents.	→		
Extraction model	<input checked="" type="radio"/> Ready	4 types trained	
Train the model to extract the data from your documents.	→		
Data standardization	<input type="radio"/> Not ready	0 types reviewed	
Map fields to new or existing data definitions.	→		
Document retention	<input checked="" type="radio"/> Ready	5 types reviewed	
Determine how long you want documents to stay in your content repository.	→		

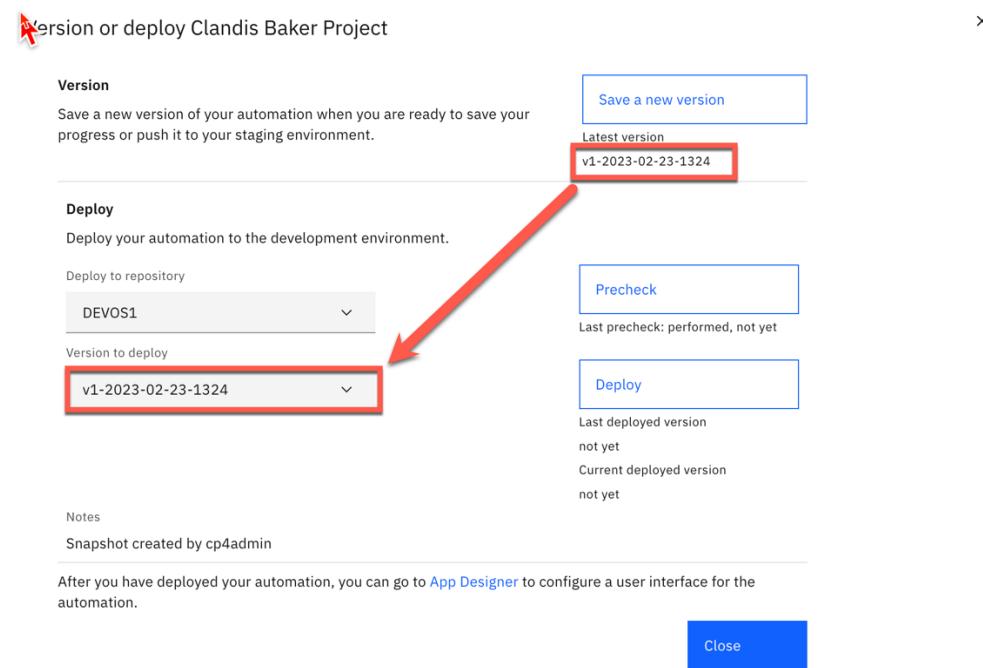
Share Version / Deploy
Last shared | 6 hours ago Latest version | not yet
Deployed | not yet

_2. Click **Save a new version**. A *Save a new version* window pops up.



_3. Click on **Create a version**

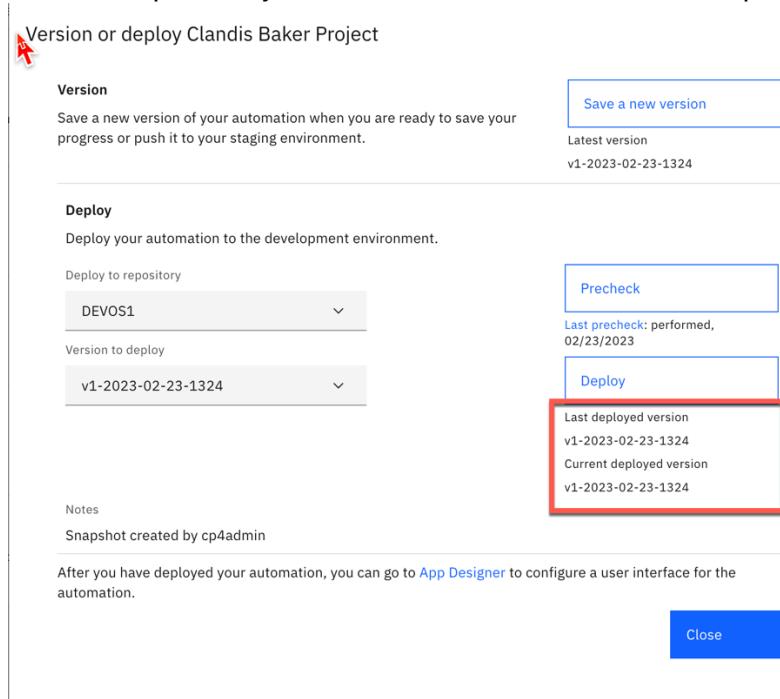
_4. Once the version is saved, you should see the version in the Version to deploy drop down list



... 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.

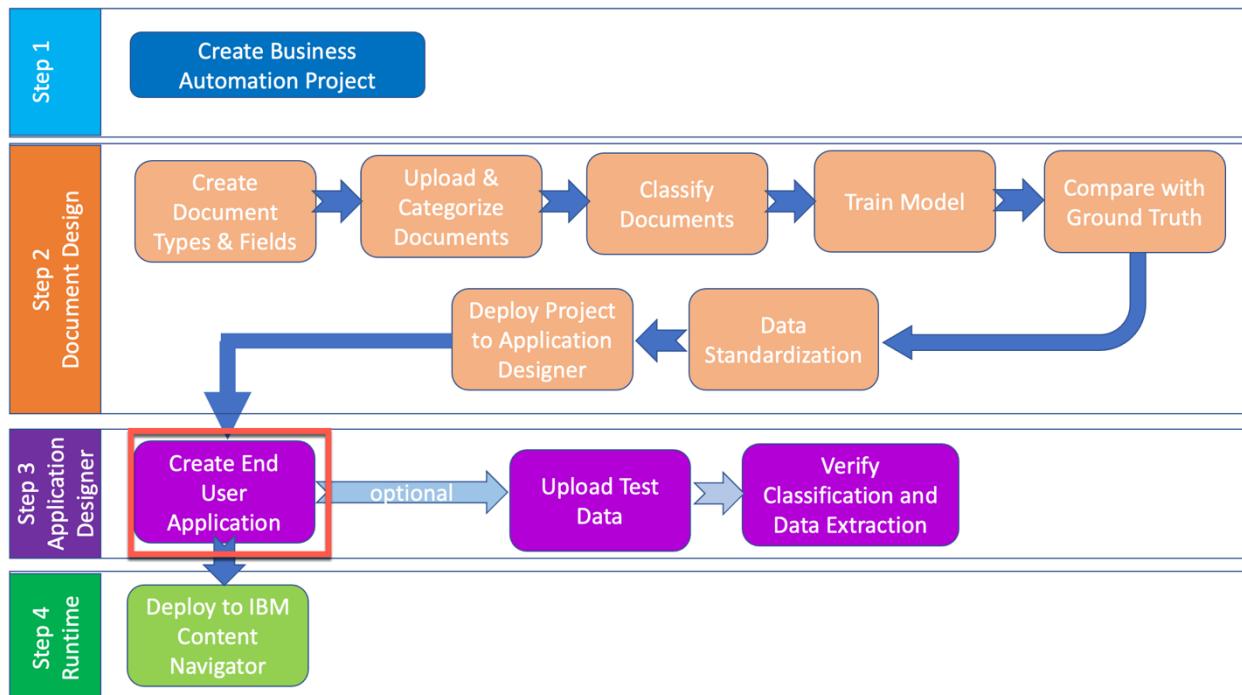
6. Click Close button

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 types	23 samples on average
Upload sample documents to define the types of documents you want the system to process.	→		
Classification model	Retrain	5 types trained	100% accuracy
Train the model to classify your documents.	→		
Extraction model	Ready	4 types trained	
Train the model to extract the data from your documents.	→		
Data standardization	Not ready	0 types reviewed	
Map fields to new or existing data definitions.	→		
Document retention	Ready	5 types reviewed	
Determine how long you want documents to stay in your content repository.	→		

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

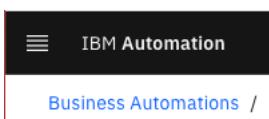
1. Batch Document Processing template – used to process batches of documents
2. Document Processing Template – used to process single documents

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.

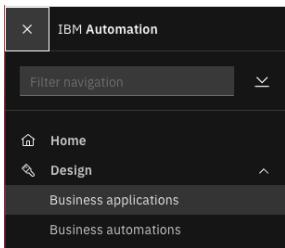
Changes to the application itself will not be in the scope of this lab.

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

Template Name	Description	Last Updated
Request Approval template	Use this template to create a service desk request.	02/20/2023
Onboarding Application template	Use this template to onboard new employees to your organization.	02/20/2023
Exception Handling template	Use this template to create a basic refund request application.	02/20/2023

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

_4. In the Create Form Template in drop down **select Batch Document Processing template (BCAT)**

Create a business application

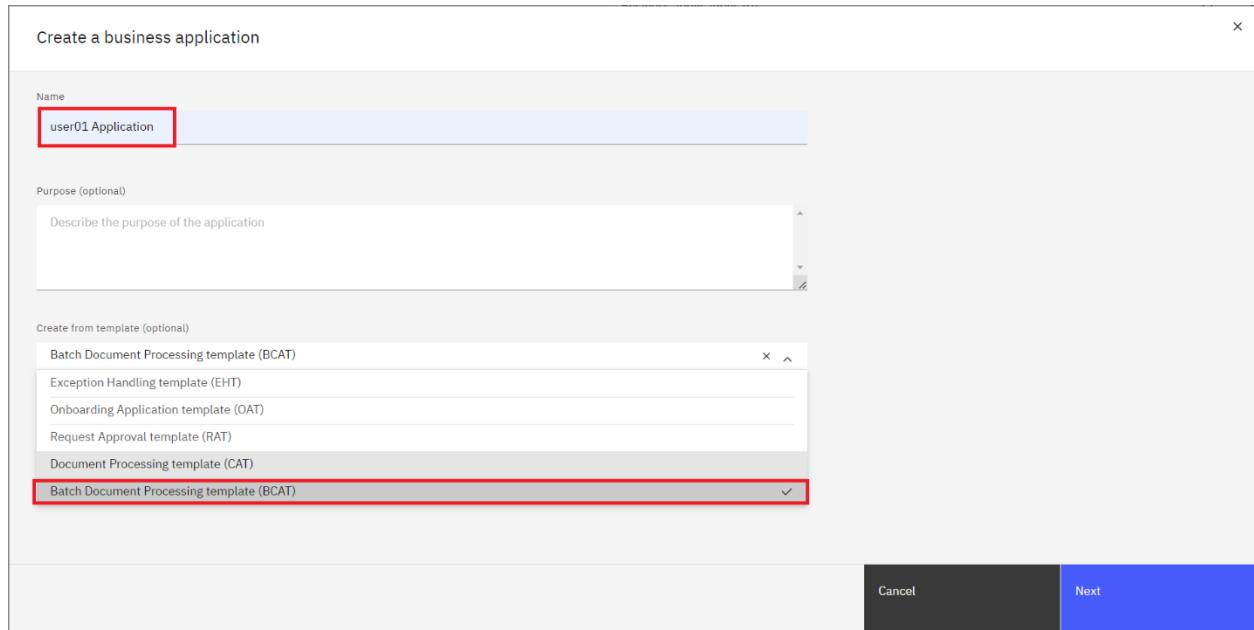
Name
user01 Application

Purpose (optional)
Describe the purpose of the application

Create from template (optional)

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

Cancel Next



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**

Create an application

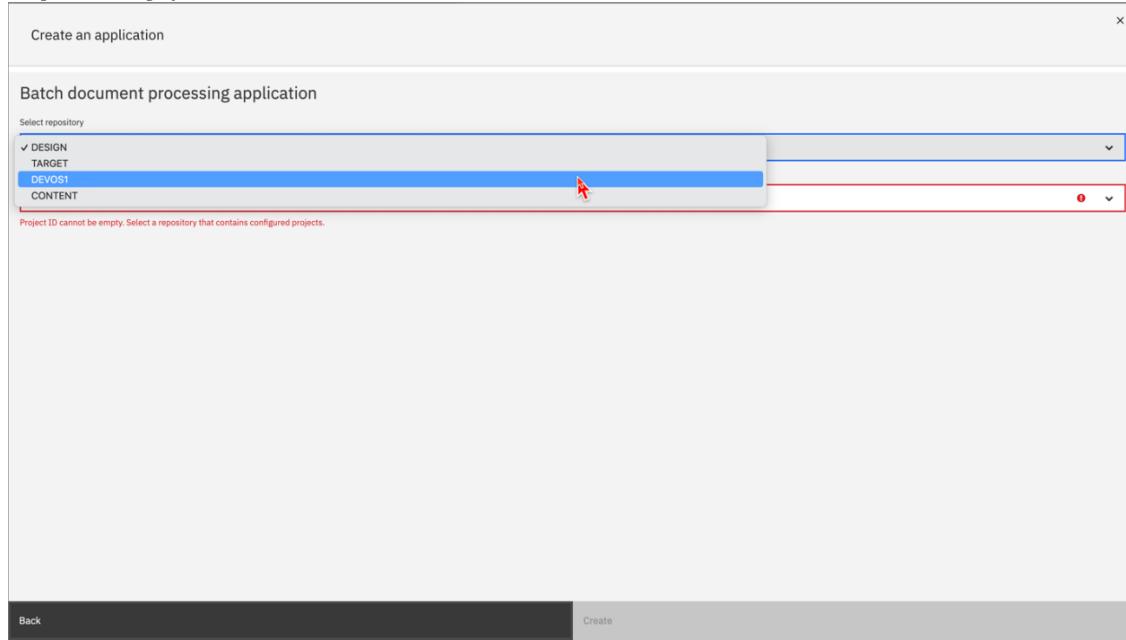
Batch document processing application

Select repository

✓ DESIGN
TARGET
DEVOS1
CONTENT

Project ID cannot be empty. Select a repository that contains configured projects.

Back Create



_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**

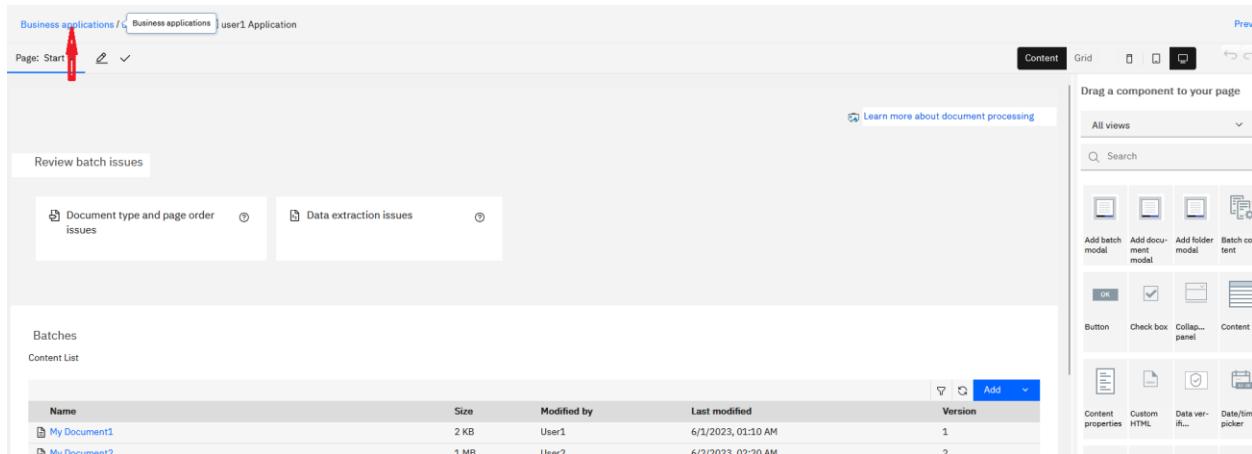
You should now be in the *Application Designer*

Name	Size	Modified by	Last modified	Version
My Document1	2 KB	User1	6/1/2023, 01:10 AM	1
My Document2	1 MB	User2	6/2/2023, 02:20 AM	2
My Document3	90 B	User3	6/3/2023, 03:30 AM	3
My Document4	1.2 MB	User4	6/4/2023, 04:40 AM	4



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

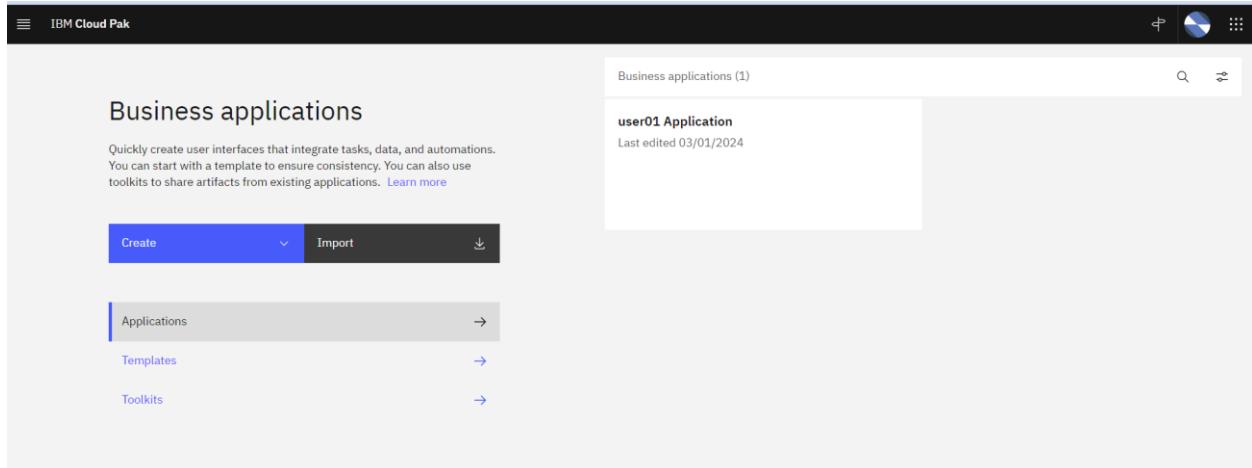


The screenshot shows the IBM Cloud Pak interface for Business applications. At the top, there is a breadcrumb navigation bar with the path: Business applications / Business applications / user1 Application. A red arrow points to the first 'Business applications' breadcrumb. Below the navigation bar, there is a header with 'Page: Start' and some icons. The main content area has a title 'Review batch issues' and two sections: 'Document type and page order issues' and 'Data extraction issues'. To the right, there is a sidebar titled 'Content Grid' with various components listed under 'Drag a component to your page' such as 'All views', 'Search', 'Batch content', 'Content list', 'Content properties', 'Custom HTML', 'Data verifier', and 'Date/time picker'. There is also a section for 'Batch models' with icons for 'Add batch model', 'Add document model', 'Add folder model', and 'Batch content'.



Note: It may take several seconds up to multiple minutes to build and display the current configuration of the interface. In case the screen does not load properly the first time, try to reload the whole browser window.

10. If you hover over any of the applications on the right, the respective box will turn grey, and a Preview and Open link will become visible. Clicking Preview would let you test the pre-configured interface. Clicking Open would open the designer for the application where you can modify the look and feel and modify its features.
Click anywhere into the grey box, but not the Preview or Open link. This brings you to the details of the application.



The screenshot shows the 'Business applications' interface. On the left, there is a sidebar with 'Create' and 'Import' buttons. Below them are links for 'Applications', 'Templates', and 'Toolkits'. On the right, there is a panel titled 'Business applications (1)' containing the details for 'user01 Application', which was last edited on 03/01/2024. The 'Applications' link in the sidebar is highlighted with a blue border, indicating it is the active section. The overall interface is clean and modern, typical of a cloud-based application management tool.

_11. From this screen if you **click** on the **3 dots** you could for example export the application or delete it

_12. Now **click** on **Preview**



Note: You may have a popup blocker turned on in your browser. Your browser will need to have this option off for the Preview.

The Preview allows you to validate the execution behavior of your application.

Previewing your application is a vital step in the creation process. You can preview your application at various points throughout your development. Maybe you want

to preview a small interaction within your application or test the entire experience of your application after you complete development.

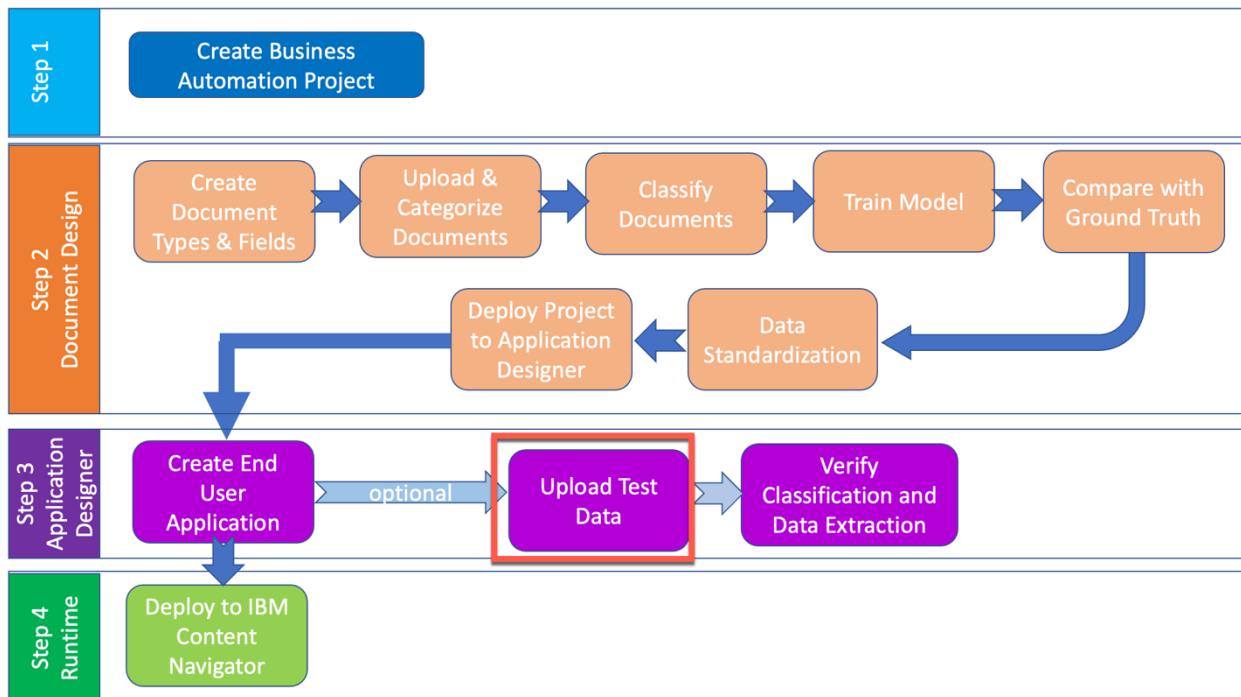
- _13. In case the Preview takes time more than about 9 minutes or throws an exception like “Unable to connect to server”, open a new browser and type the URL -> <https://cpd-ibm-cp4ba.apps.ocp.ibm.edu/ae-pbk/v2/applications>

This should look like below –

```
JSON Raw Data Headers
Save Copy Collapse All Expand All Filter JSON
apps:
  0:
    name: "ttt"
    uniqueName: "ttt(TT)"
    id: "28c49208-7f45-45f5-876a-71badff40ae46"
    url: "https://cpd-ibm-cp4ba.apps.ocp.ibm.edu/ae-pbk/ttt(TT)" [red box]
    isPublicApp: false
    iconUrl: ""
    display: true
    description: ""
    lastModified: "2024-04-03T22:41:19.520Z"
    versionName: ""
    translations: []
    exposedAs: "App"
    lastPublished: "2024-04-03T22:46:22.448Z"
    teams: []
  1:
    name: "Client Onboarding Document Upload"
    uniqueName: "Client Onboarding Document Upload(CODU)"
    id: "2b16c93c-2ad3-41a4-b65d-cfceac21e1f6"
    url: "https://cpd-ibm-cp4ba.apps.ocp.ibm.edu/ae-pbk/Client Onboarding Document Upload(CODU)"
    isPublicApp: true
    publicAppUrl: "https://cpd-ibm-cp4ba.apps.ocp.ibm.edu/ae-pbk/public-app/Client Onboarding Document Upload(CODU)"
    iconUrl: ""
```

Here you can observe the ADP application that you created. When you refer to the above snapshot, you can see that there is an application called “ttt” and notice the URL. Copy the URL and paste it into a new browser window.

11.2 Upload documents for processing

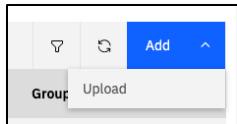


- _1. Below pasted snapshot is the preview of the application. Normally, this preview should work in “incognito mode” of Chrome or “In Private mode” window of an Edge browser. Additionally, any popup blocker must be disabled or configured to allow open the pop-up window. You should be in the default application user interface for ADP. It opens a new tab/window in your browser.

The screenshot shows the IBM Content Navigator interface. At the top, there's a navigation bar with icons for file operations like Open, Save, Print, and a search bar. The main area has two sections: "Review batch issues" and "Batches". The "Review batch issues" section shows two boxes: "Document type and page order issues" (0 batches) and "Data extraction issues" (0 batches). The "Batches" section has a header with columns: Name, Files, Priority, Status, Added on, Added by, Group, and Location. A search icon is at the top of the list. Below it, a message says "No items found." There are also "Add" and "Import" buttons.

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. The example below shows you how you would 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

Name	Files	Priority	Status	Added on	Added by	Group	Location
Batch01	5	High		3 of 5 files processed	02/23/2023, 10:49 AM	cp4admin	

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

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

Name	Files	Priority	Status	Added on	Added by	Group	Location
Batch01	5	High		Documents uploaded	02/23/2023, 10:49 AM	cp4admin	

_8. Click Submit

In the screen shot below, you see the status of the batch job is marked as having Document issues. Matching with that 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 get to the respective batches

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 **red checkmark** for documents that have a **low confidence** document type
- a **red exclamation mark** for documents that **could not be classified**

ITEM #	DESCRIPTION	UNIT PRICE	LINE TOTAL
01	Whole Chicken	£1.00	£4.00
02	One Day Old Chick	£1.00	£1.00

- _3. Most of the document types are correct but it looks like a Purchase order (PO) got mixed into our batch. **Click** on the **Trash can** to delete it from the batch and **select OK** to finally delete it.

The screenshot shows the 'Batch01' interface. On the left, a list of documents is displayed under 'Documents (5)'. One document, 'B_PO_5.pdf', has its 'Document type' set to 'Undefined' and features a trash can icon. On the right, a detailed view of a 'PURCHASE ORDER' document is shown. The document header includes 'RUBE'S Meat Co.', 'P.O. No.: 71230', 'DATE: 09 March 2020', and 'CUSTOMER ID: 447320'. The 'SHIP TO:' section lists 'Chicken Run Ranch' with address '24 Old Quay Lane, Nelson Village NE23 6DD, UK' and contact '078-2064-8486'. The 'SHIPPING METHOD' table shows 'AIR' selected. The 'DELIVERY DATE' table shows '29 March 2020'. The main body of the document contains two rows of items:

QTY	ITEM #	DESCRIPTION	JOB	UNIT PRICE	LINE TOTAL
230 PCS	01	Whole Chicken		£1.50	£345.00
150 Packs	02	One Day Old Chick		£1.05	£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.

[Batch Document Processing Application](#) /
Document type and page order issues

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

A screenshot of a user interface tile titled "Data extraction issues". Below the title, it says "1 batches". There is a small question mark icon in the top right corner.

_2. Click on <your Batch name> to open

A screenshot of a list titled "Name". It contains one item: "Batch 1".

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

A screenshot of a table showing a list of documents. The table has columns: Name, Issues, Status, Modified on, and Modified by. One row shows a document named "BAD_FW2_1000_0003F.pdf" with 1 issue, status "Data issues", modified on 03/04/2023, and modified by "cp4admin". At the bottom, there are navigation links for "Submit" and "Items per page: 100 1-4 of 4 items".

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

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

The screenshot shows a document processing interface with a PDF viewer on the left and an "Extracted data" panel on the right.

Document type: Wage and Tax

Extracted data

- Federal Income Tax Withheld:** 9000.00
- Employee Social Security Number:** * (Validation error)
- 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, Employee's address and ZIP code Stella K. James 343 Twisting Way Red Beach, CA 90354

Take a moment to discover the image viewer features.

Image viewer features at top:

The screenshot shows a document processing interface with a PDF viewer on the left and an "Extracted data" panel on the right.

Document type: Wage and Tax

Extracted data

- Employee Name and Address:** * (Validation warning)
- Name:** (none)
- Email:** (none)
- Phone:** (none)
- Postal mail address:**
 - Building number:** 457
 - Street name:** Chelsea Place
 - Unit:** (none)

- Rotate image
- Visual effect adjustment
- Invert

Image viewer features at bottom:

The screenshot shows a document processing application window. At the top, it displays the file name "BAD_FW2_1000_0003F.pdf" and the document type "Wage and Tax". On the right side, there is an "Extracted data" panel with sections for "Federal Income Tax Withheld" (value: 9000.00) and "Employee Social Security Number" (value: (none)). Below this is a "Similar fields" section. A red box highlights the "Extracted data" header and the "All Fields" dropdown.

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

- Page and thumbnail's view
- Fit to window
- Zoom and Magnify

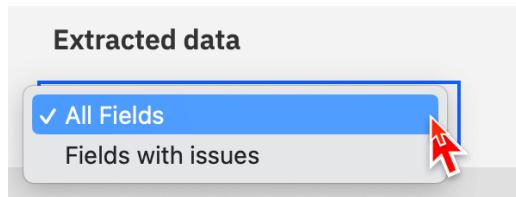
Field features

This screenshot shows the same document processing application window, but the "Extracted data" panel is no longer visible. The rest of the interface remains the same, including the document preview, navigation controls, and the "Save changes" button.

- Show all fields.
- Show fields with issues.

Also note that fields that do 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**

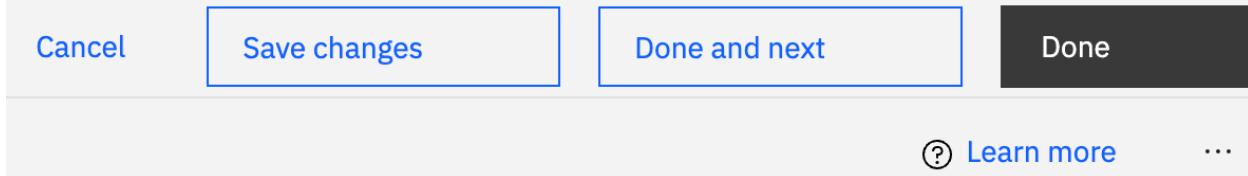


The Employee Social Security Number is a mandatory field. For purposes of this lab it was changed to "Bad SSN". Since you did not make that phrase an alias ADP was not able to pick it up.

6. Click on Employee Social Security Number and with your mouse select the SSN under "Bad SSN"

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



- _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 Export/Import Project (Optional)

If you would like to save your project and perhaps use it later, you can perform the steps in this chapter.

From the Business Automations

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

The screenshot shows the IBM Cloud Pak Business Automations interface. On the left, there's a sidebar with 'Business automations' and a 'Create' button. In the center, a list titled 'Document processing automations (3)' shows an item named 'User01_CEB' last edited on 09/20/2023. At the bottom, there are tabs for 'Published automation services' (with 'Document processing' selected) and 'External'.

_2. Select <your project name>. Click open

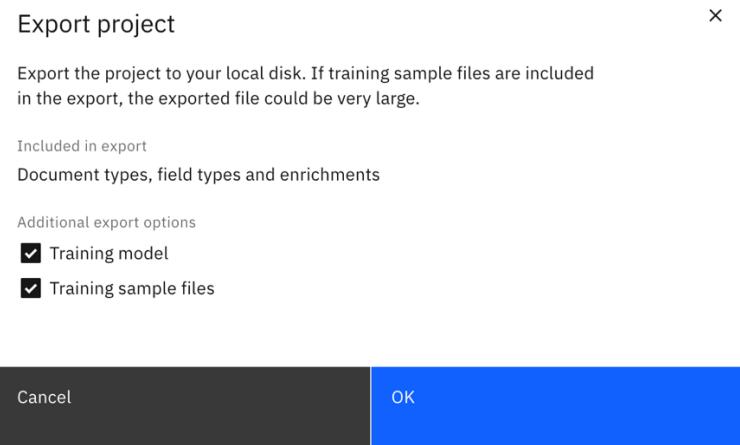
The screenshot shows the same interface as above, but the 'User01_CEB' project is now highlighted with a gray background. A blue 'Open' button is visible at the bottom right of its card.

_3. From the main screen **select the Configure tab**

The screenshot shows the 'Configure' tab selected in the top navigation bar. The left sidebar has 'Import / Export ontology' expanded, showing 'Document processing', 'Language settings', 'Git server configuration', and 'Webhook configuration'. The right side shows 'Export project' and 'Import project' buttons.

_4. Select Export Project

_5. On Export Project window **check Training modul** and **Training sample files**



_6. Click on OK

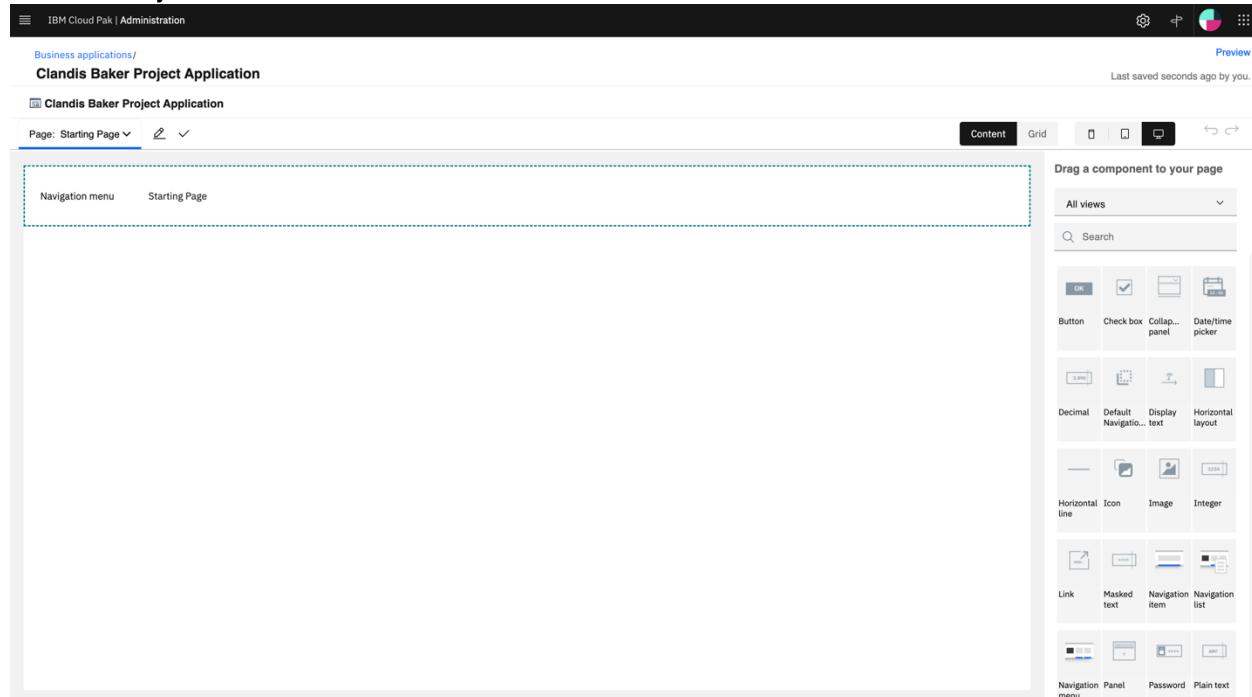
_7. A project-export-<date-time>.zip will be download via browser to local machine.

You have successfully completed the Automation Document Processing lab.
Congratulations and well done!

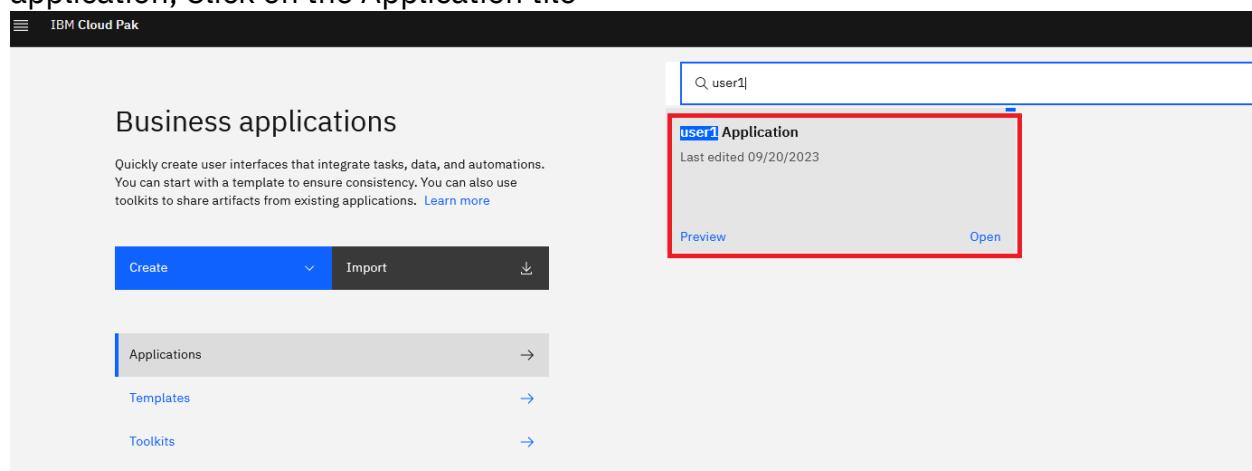
Appendix A - Troubleshooting

Application Blank

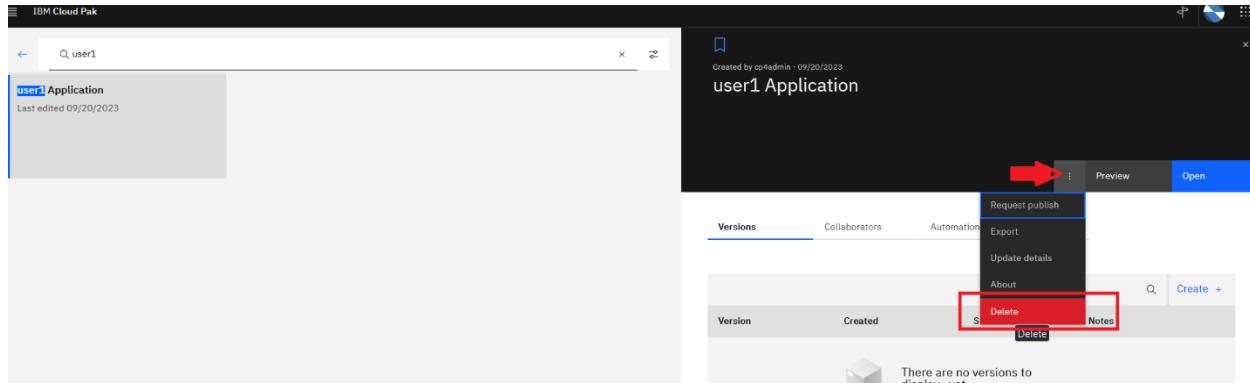
During creating of Business Application setup, sometimes on first time after project has been deployed. The Starter page remains blank or shows the loading animation indefinitely.



First try to reload the whole editor page and wait for the UI to be loaded. If this remains unsuccessful, delete the application and try again. To delete the application, Click on the Application tile

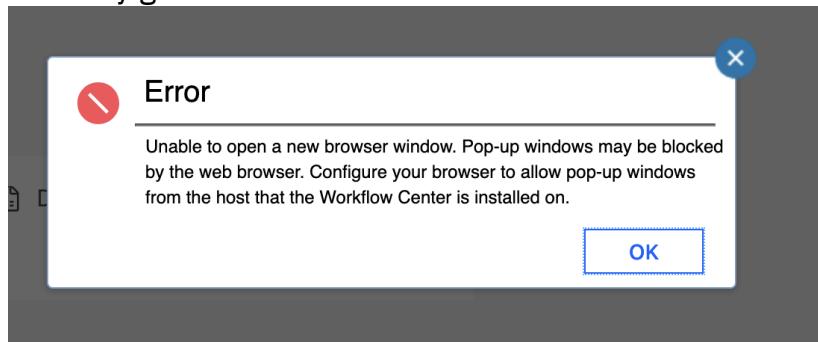


Then click on the 3 dots and select Delete.



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>.