

## Google Cloud Skills Boost for Partners

[Main menu](#)

## Build Custom Processors with Document AI

Course · 6 hours 66%  
15 minutes complete

## Build Custom Processors with Document AI

Optical Character Recognition (OCR) with Document AI (Python)

Form Parsing with Document AI (Python)

Using Specialized Processors with Document AI (Python)

Uptraining with Document AI Workbench

Custom Document Extraction with Document AI Workbench

Course &gt; Build Custom Processors with Document AI &gt;

Quick tip: Review the prerequisites before you run the lab

[End Lab](#)

00:36:42

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked.  
[Learn more.](#)[Open Google Cloud console](#)

Username

student-00-d6dc143dd07cf

Password

d7tz4xzaFrNJ

Project ID

qwiklabs-gcp-04-3ead50e1

# Custom Document Extraction with Document AI Workbench

 Lab
 1 hour
 No cost
 Intermediate

★★★★★

[This lab may incorporate AI tools to support your learning.](#)

GSP1142

- Document AI API
- Task 2. Create a processor 100/100
- Task 3. Define processor fields
- Task 4. Upload a sample document
- Task 5. Label a document
- Task 6. Build processor version using foundation model
- Task 7. Use generative AI to auto-label documents
- Task 8. Import prelabeled training documents
- Task 9. Train the processor
- Congratulations!

[Previous](#)[Next >](#)

## Overview

Document AI is a document understanding solution that takes unstructured data (e.g. documents, emails, invoices, forms, etc.) and makes the data easier to understand, analyze, and consume. The API provides structure through content classification, entity extraction, advanced searching, and more. With Document AI Workbench, you can achieve higher document processing accuracy by creating fully customized models using your own training data.

You can create Custom Document Extractors (CDE) that are specifically suited to your documents, and trained and evaluated with your data. This processor identifies and extracts entities from your documents. You can then use this trained processor on additional documents. You typically would use a CDE on documents that are all of one type, such as your institution's enrollment forms.

CUSTOM DOCUMENT EXTRACTOR THAT PROCESSES W-2 (US TAX FORM) DOCUMENTS. MOST OF THE DOCUMENT PREPARATION WORK HAS BEEN DONE SO THAT YOU CAN FOCUS ON THE OTHER MECHANICS OF CREATING A CDE.

## Objectives

In this lab, you will learn how to perform the following tasks:

- Create a Custom Document Extractor in Document AI Workbench
- Define and create the processor schema
- Import documents
- Annotate documents manually in Document AI Workbench
- Use generative AI to auto-label documents
- Kick off a training job for the processor

## Setup and requirements

Before you click the Start Lab button

Read these instructions. Labs are timed and you cannot pause them. The timer, which

Read these instructions; labs are timed and you cannot pause them. The timer, which starts when you click **Start Lab**, shows how long Google Cloud resources are made available to you.

This hands-on lab lets you do the lab activities in a real cloud environment, not in a simulation or demo environment. It does so by giving you new, temporary credentials you use to sign in and access Google Cloud for the duration of the lab.

To complete this lab, you need:

- Access to a standard internet browser (Chrome browser recommended).

This prevents conflicts between your personal account and the student account, which may cause extra charges incurred to your personal account.

- Time to complete the lab—remember, once you start, you cannot pause a lab.

**Note:** Use only the student account for this lab. If you use a different Google Cloud account, you may incur charges to that account.

## How to start your lab and sign in to the Google Cloud console

1. Click the **Start Lab** button. If you need to pay for the lab, a dialog opens for you to select your payment method. On the left is the Lab Details pane with the following:

- The Open Google Cloud console button
- Time remaining

- Other information, if needed, to step through this lab

2. Click **Open Google Cloud console** (or right-click and select **Open Link in Incognito Window** if you are running the Chrome browser).

The lab spins up resources, and then opens another tab that shows the Sign in page.

*Tip:* Arrange the tabs in separate windows, side-by-side.

**Note:** If you see the **Choose an account** dialog, click **Use Another Account**.

3. If necessary, copy the **Username** below and paste it into the **Sign in** dialog.

student-00-d6dc143dd07c@qwiklabs.net



You can also find the Username in the Lab Details pane.

5. Copy the **Password** below and paste it into the **Welcome** dialog.

d7tz4xzaFrNj



You can also find the Password in the Lab Details pane.

6. Click **Next**.

**Important:** You must use the credentials the lab provides you. Do not use your Google Cloud account credentials.

**Note:** Using your own Google Cloud account for this lab may incur extra charges.

7. Click through the subsequent pages:

- Accept the terms and conditions.

a temporary account).

- Do not sign up for free trials.

After a few moments, the Google Cloud console opens in this tab.

**Note:** To access Google Cloud products and services, click the **Navigation menu** or type the service or product name in the **Search** field.



DASHBOARD ACTIVITY RECOMMENDATIONS

## Activate Cloud Shell

Cloud Shell is a virtual machine that is loaded with development tools. It offers a persistent 5GB home directory and runs on the Google Cloud. Cloud Shell provides command-line access to your Google Cloud resources.

2. Click through the following windows:

- Continue through the Cloud Shell information window.
- Authorize Cloud Shell to use your credentials to make Google Cloud API calls.

When you are connected, you are already authenticated, and the project is set to your **Project\_ID**, `qwiklabs-gcp-04-3ead50e8fb32`. The output contains a line that declares the **Project\_ID** for this session:

```
Your Cloud Platform project in this session is set to qwiklabs-gcp-04-3ead50e8fb32.
```

`gcloud` is the command-line tool for Google Cloud. It comes pre-installed on Cloud Shell and supports tab-completion.

3. (Optional) You can list the active account name with this command:

4. Click **Authorize**.

**Output:**

```
ACTIVE: *
ACCOUNT: student-00-d6dc143dd07c@qwiklabs.net

To set the active account, run:
$ gcloud config set account `ACCOUNT`
```

5. (Optional) You can list the project ID with this command:

```
gcloud config list project
```

**Output:**

```
[core]
```

**Note:** For full documentation of `gcloud`, in Google Cloud, refer to the [gcloud CLI overview guide](#).

## Task 1. Enable the Document AI API

Before you can begin using Document AI, you must enable the API.

1. In Cloud Shell, run the following command to enable the API for Document AI.

```
gcloud services enable documentai.googleapis.com
```

You should see something like this:

```
Operation "operations/..." finished successfully.
```

2. Run the following command to install the Python client libraries for Document AI.

```
pip3 install --upgrade google-cloud-documentai
```

You should see something like this:

```
...
Installing collected packages: google-cloud-documentai
Successfully installed google-cloud-documentai-2.15.0
```

Now, you're ready to use the Document AI API!

Enable the Document AI API



Check my progress

## Task 2. Create a processor

You must first create an Custom Document Extractor processor to use for this lab.

You must first create a Form Parser processor instance to use in the Document AI Platform for this tutorial.

1. From the Navigation menu select **View all products**. Under **Artificial Intelligence**, select **Document AI**.

The screenshot shows the Google Cloud Document AI Overview page. It includes a 'Get started with Document AI' section, a 'How it works' section with three numbered steps (Create a processor, Get structured data, Use), and a 'Resources' section with links to 'What's new in Document AI', 'Documentation', 'Technical resources', and 'Training information'.

2. Click **Create Custom Processor**.

3. Inside the **Custom Extractor** box, click **Create Processor**.

4. Give it the name **lab-custom-extractor** and select the region **US (United States)** on the list.

5. Click **Create** to create your processor.

Click **Check my progress** to verify the objective.

Assessment Completed!

## Task 3. Define processor fields

You are now on the **Processor overview** page of the processor you just created.

The screenshot shows the Google Cloud Document AI Processor details page for 'lab-custom-extractor'. It includes sections for 'Processor details' (Overview, Get started, Build, Evaluate & test), 'Using your processor' (Use out-of-the-box, Customize), and 'Basic information'.

Name	lab-custom-extractor
ID	fddcd39e039978c7b1
Status	<input checked="" type="checkbox"/> Enabled
Processor Type	Custom Extractor
Created	Oct 12, 2023, 11:47:23 AM
Encryption Type	Google-managed
Region	us

You can specify the fields you want the processor to extract and begin labeling documents.

1. Click on the **Get started** tab. The **Fields** menu appears.
2. Click **Create New Field**.
3. Enter the name for the field. Select the **Data type** and the **Occurrence**. Click **Create**. Refer to [Define processor schema](#) for detailed instructions on creating and editing a schema.
4. Create each of the following labels for the processor schema.

Name	Data Type	Occurrence
employees_social_security_number	Number	Required multiple
employer_identification_number	Number	Required multiple
employers_name_address_and_zip_code	Address	Required multiple
federal_income_tax withheld	Money	Required multiple
social_security_tax withheld	Money	Required multiple
social_security_wages	Money	Required multiple
wages_tips_other_compensation	Money	Required multiple

You can also create and use [other types of labels](#) in your processor schema, such as checkboxes and tabular entities. For example, the W-2 forms contain a **Statutory employee**, **Retirement plan**, and **Third party sick pay** check boxes that you could also add to the schema.

Field	Type	Occurrence	Extract	Required	Multiple
employer_identification_number	Number	Required multiple	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
employers_name_address_and_zip_code	Address	Required multiple	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
federal_income_tax withheld	Money	Required multiple	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
social_security_tax withheld	Money	Required multiple	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
social_security_wages	Money	Required multiple	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
wages_tips_other_compensation	Money	Required multiple	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

Click **Check my progress** to verify the objective.

Create Labels

Check my progress

Assessment Completed!

## Task 4. Upload a sample document

1. Click **Upload Sample Document**.
2. In the sidebar, click **Import documents from Google Cloud Storage**.
3. For this example, enter this bucket name in Source path. This links directly to one document.

Import

cloud-samples-data/documentai/Custom/W2/PDF/W2\_XL\_input\_clean\_2950.pdf

4. Click **Import**.

You are redirected to the labeling console.

The process of selecting text in a document, and applying labels is known as annotation.

- When you're at the labeling console, notice that many of the labels are already populated.

**Note:** Your results might look slightly different than the sample image.

- To use the suggested labels, hold the pointer over each label in the side panel, and click on the check mark to confirm the label is correct. You can edit the values if they do not match the document text.
- In this example, the values at the bottom of the document were not identified automatically, so you will need to label them manually.
- Use the **Bounding box** tool by default, or the **Select text** tool for multi-line values, to select the content and apply the label.

**Note:** The **Select text** tool does not work for all text values, so use the **Bounding box** if appropriate. You can also select non-text fields such as checkboxes using the **Bounding box** tool.

- Review the detected text values to ensure that they reflect the correct text from the document.

The labeled W-2 document should look like this when complete:

--	--

7. If needed, you can click **Create New Field** to add a new field to the schema from this page.

8. Click **Mark as Labeled** when you have finished annotating the document.

You are redirected to the **Get started** tab.

## Task 6. Build processor version using foundation model

After labeling a single document, you can create a processor version using the pretrained foundation model to extract entities.

The screenshot shows the 'my-custom-document-extractor' processor details page. The 'Build' tab is active. The 'Dataset overview' section shows 1 total document, 1 labeled document, 0 auto-labeled documents, and 0 unlabeled documents. Below this, there are two main options: 'Call foundation model' (which creates a version with zero training using Google's foundation model and fields you created) and 'Train a custom model' (which trains a custom model from scratch using a labeled dataset and the fields you created). At the bottom, there are 'VIEW FULL REQUIREMENTS', 'CREATE NEW VERSION', and 'CREATE NEW VERSION' buttons.

2. Under **Call foundation model**, click **Create New Version**.

3. Enter a name for your processor version, such as `w2-foundation-model`.

**Note:** Once you create a processor version, you cannot delete fields you have created. You can disable them on the fields page if you no longer need them.

5. Optional: Click on the **Deploy & Use** tab. On this page, you can view the available processor versions and the deployment status of the new version.

You test and evaluate this version later in the lab.

Click **Check my progress** to verify the objective.

The dialog box has a green checkmark icon and a 'Check my progress' button. Below the button, the message 'Assessment Completed!' is displayed.

## Task 7. Use generative AI to auto-label documents

The foundation model can accurately extract fields for a variety of document types, but you can also provide additional training data to improve the accuracy of the model for specific document structures.

Document AI Workbench uses the label names you define and previous annotations to make it quicker and easier to label documents at scale with [auto-labeling](#).

1. Go to the **Build** page.

2. Click **Import Documents**.

3. In the sidebar, click **Import documents from Google Cloud Storage**.
4. Enter this bucket name in **Source path**. This contains unlabeled W-2 PDF files.



5. From the **Data split** list, select **Auto-split**. This automatically splits the documents to have 80% in the training set and 20% in the test set.
6. In the **Auto-labeling** section, select the **Import with auto-labeling** checkbox.
7. Select the foundation model processor version you just created to label the documents.
8. Click **Import** and wait for the documents to import. You can leave this page and return later.
9. You must verify the auto-labeled documents before you can use them for training or testing. Click **Start Labeling** to view the auto-labeled documents.
10. To use the suggested labels, hold the pointer over each annotation, and click on the check mark to confirm the label is correct. You can edit the values if they do not match the document text.
11. Click **Mark as Labeled** when you have finished annotating the document.
12. Repeat for each auto-labeled document. For this tutorial, you can skip any

## Task 8. Import prelabeled training documents

In this lab, you are provided with prelabeled data. If working on your own project, you have to determine how to label your data. Refer to [Labeling options](#) for more details. In general, more training data produces higher accuracy.

1. Go to the **Build** page.
2. Click **Import Documents**.
3. In the sidebar, click **Import documents from Google Cloud Storage**.
4. Enter the following path in **Source path**. This bucket contains prelabeled documents in the **Document.JSON** format



5. From the **Data split** list, select **Auto-split**. This automatically splits the documents to have 80% in the training set, and 20% in the test set. Leave **Import with auto-labeling** unchecked.
6. Click **Import**. Import takes several minutes.
7. (Optional) From the **Build** page, you can access the **Manage Dataset** console to view and edit all documents and labels in the dataset.

## Task 9. Train the processor

Now that you have sufficient training and test data, you can train the processor.

1. Under **Train a custom model**, click **Create New Version**.  
If **Create New Version** cannot be clicked, click on **View Full Requirements** for information about the dataset requirements.
2. In the **Version name** field, enter a name for this processor version, such as `w2-custom-model`.
3. (Optional) Click **View Label Stats** to find information about the document labels. That can help determine your coverage. Click **Close** to return to the training setup.

4. Under **Model training method**, select **Model based**.

5. Click **Start training**.

6. (Optional) Click on the **Deploy & Use** tab. On this page, you can view the available processor versions and the training status of the new version.

The screenshot shows the 'lab-custom-extractor' project in the Google Cloud Platform. The 'Deploy & Use' tab is selected. A table lists four versions:

Version ID	Created	Status	Name	Type
v2-custom-model	Oct 12, 2023, 12:38:58 PM	Training...	v2-custom-model	Generative AI
w2-foundation-model	Oct 12, 2023, 12:39:21 PM	Deployed	w2-foundation-model	Generative AI
pretrained-foundation-model-v1.0-2023-08-22	Aug 21, 2023, 5:00:00 PM	Deployed	Foundation Model	Generative AI

Click **Check my progress** to verify the objective.

The screenshot shows the 'Check my progress' step. It includes a green checkmark icon and the message 'Assessment Completed!'. Below this, it says 'Great! You have now started training your first Custom Document AI Processor. Since the training job will take around a few hours, the lab will end here. If you are interested in'

## Congratulations!

Congratulations, in this lab you've successfully used Document AI to create a Custom Document Extraction processor, import a dataset, and label example documents. You can now use this processor to parse documents in this format just as you would for any Specialized Processor. You can also use this processor to label new documents using auto-labeling as well as use the Document AI Workbench to manage your training data and training jobs.

## Next steps / Learn more

Check out the following resources to learn more about Document AI and the Python

- [The Future of Documents - YouTube Playlist](#)
- [Document AI Documentation](#)
- [Document AI Python Client Library](#)
- [Document AI Samples](#)

## Google Cloud training and certification

...helps you make the most of Google Cloud technologies. [Our classes](#) include technical skills and best practices to help you get up to speed quickly and continue your learning journey. We offer fundamental to advanced level training, with on-demand, live, and virtual options to suit your busy schedule. [Certifications](#) help you validate and prove your skill and expertise in Google Cloud technologies.

**Manual Last Updated March 17, 2025**

**Lab Last Tested March 17, 2025**

Copyright 2025 Google LLC. All rights reserved. Google and the Google logo are trademarks of Google LLC. All other company and product names may be trademarks of the respective companies with which they are associated.