Step 0: Create a new project
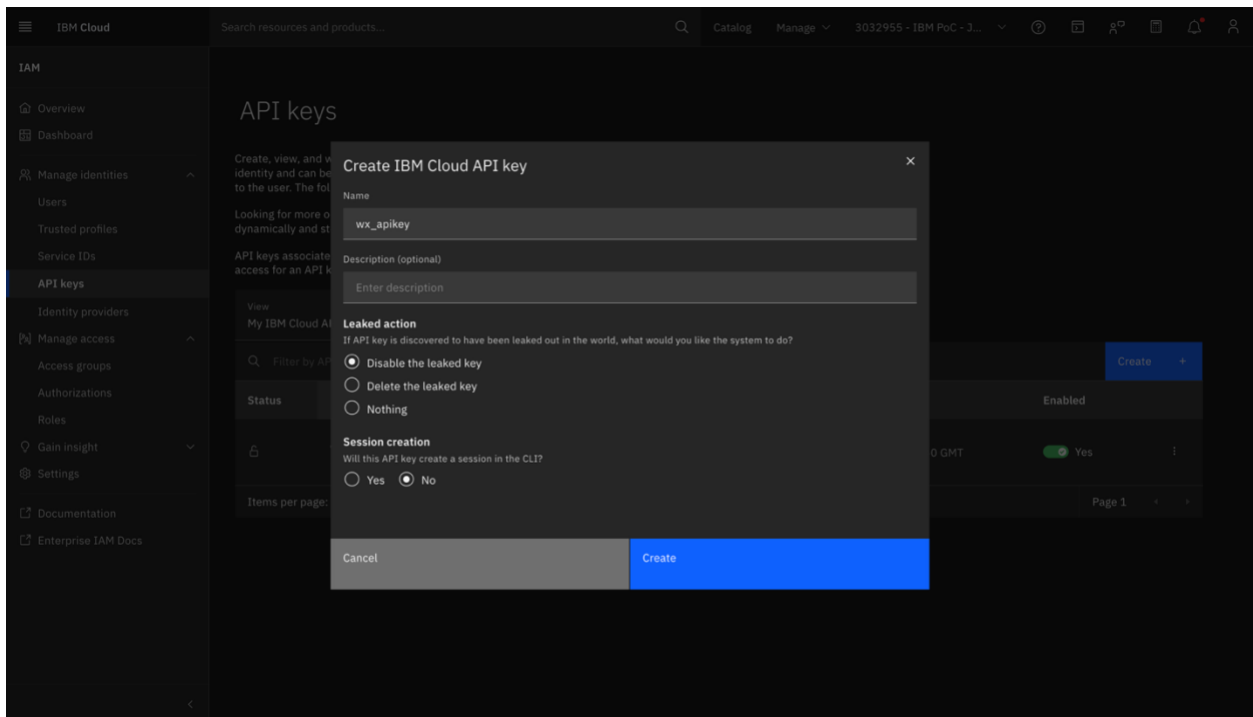
1. Open your browser and go to IBM Cloud Login: https://cloud.ibm.com/login
2. Sign in with your IBM Cloud credentials.
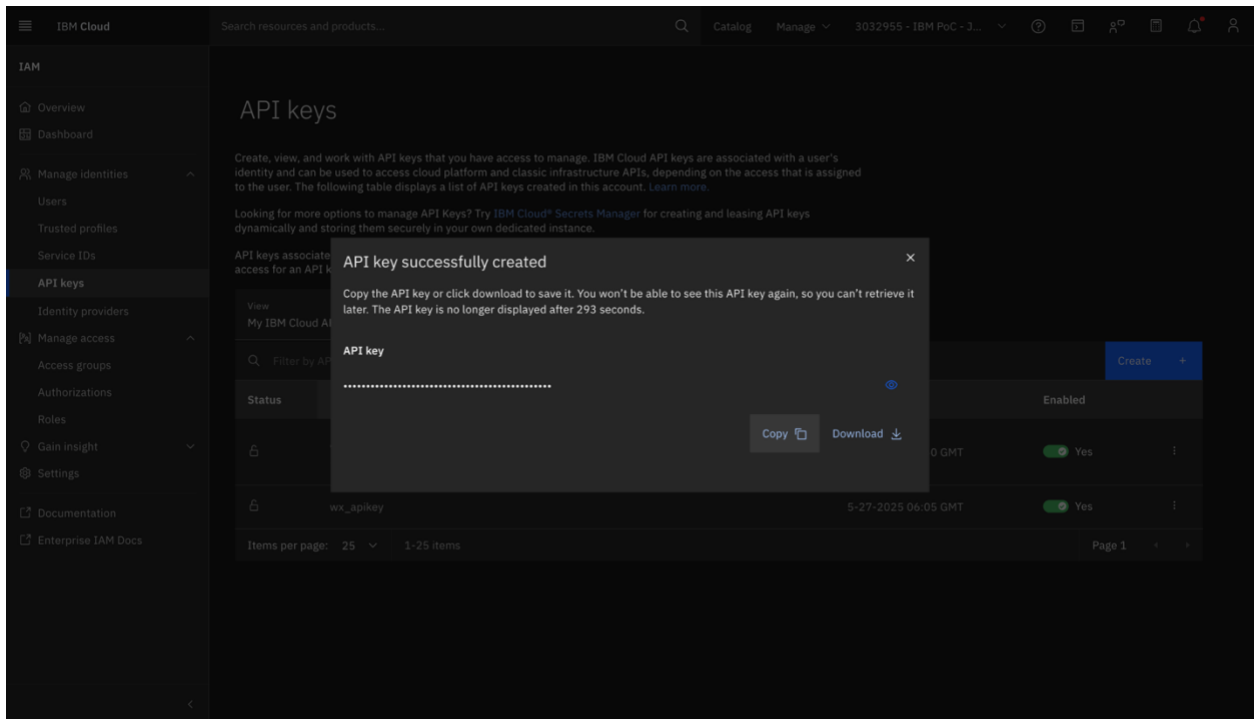3. Expand Manage → click Access (IAM).
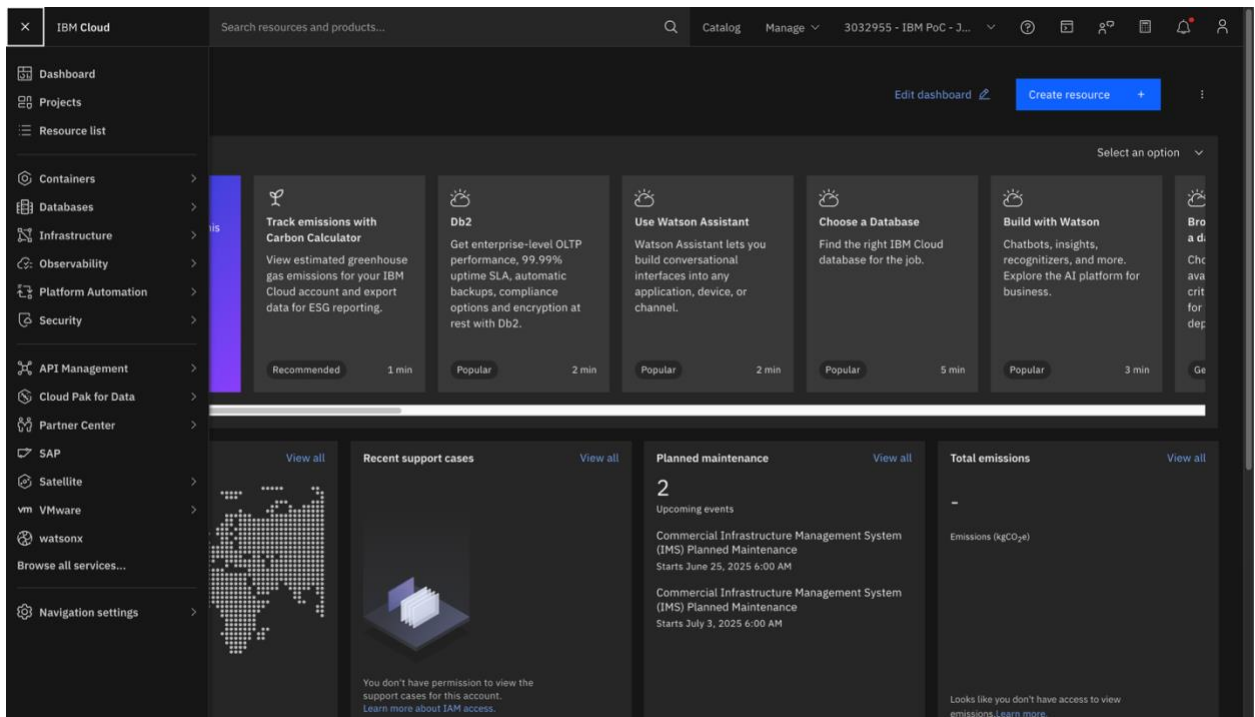


4. Select API Keys.

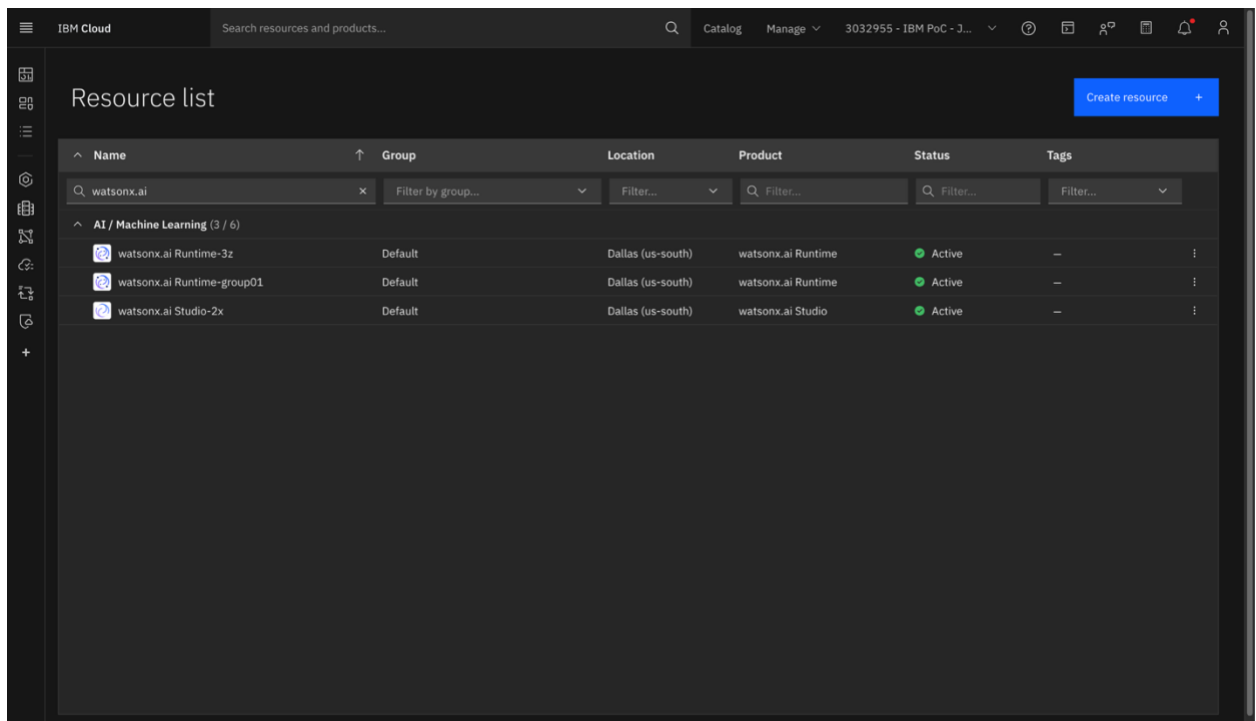5. Click Create and enter the name of key: wx_apikey.



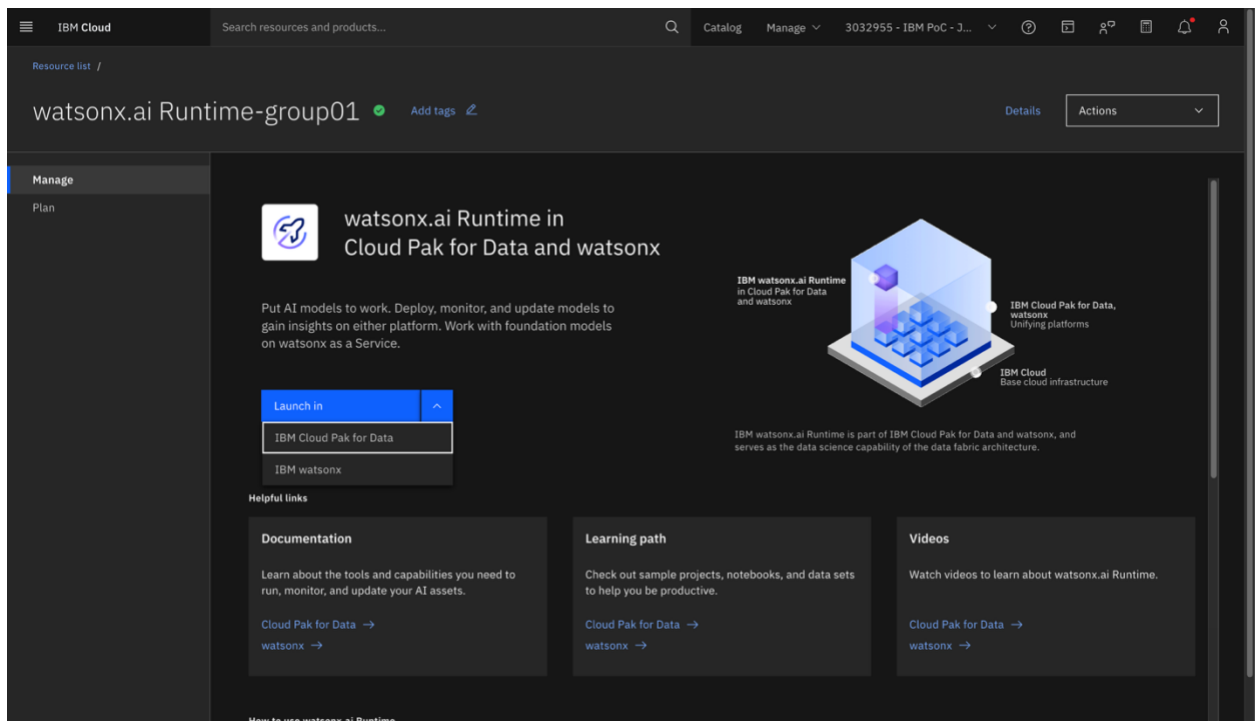6. Copy the displayed API key and save it **securely**.

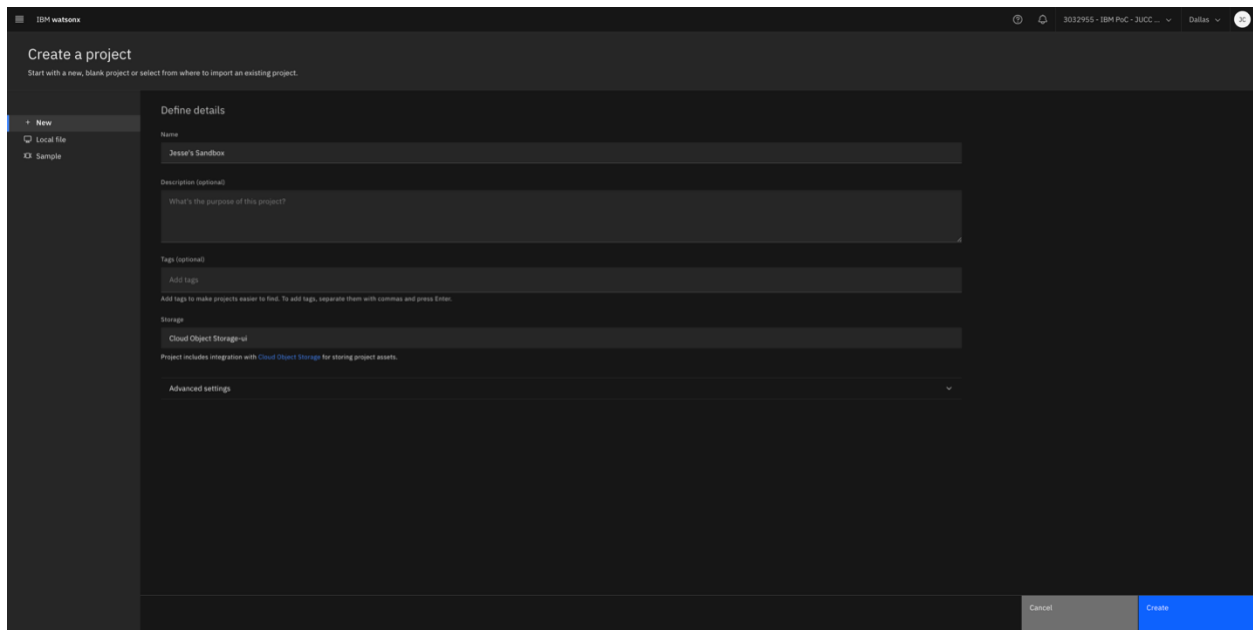7. In the left-hand sidebar, click Resource list.



8. In the search box, type watsonx.ai, then from the results select your assigned service watsonx.ai group.

9. Click Launch in → **IBM watsonx** to open the watsonx console.

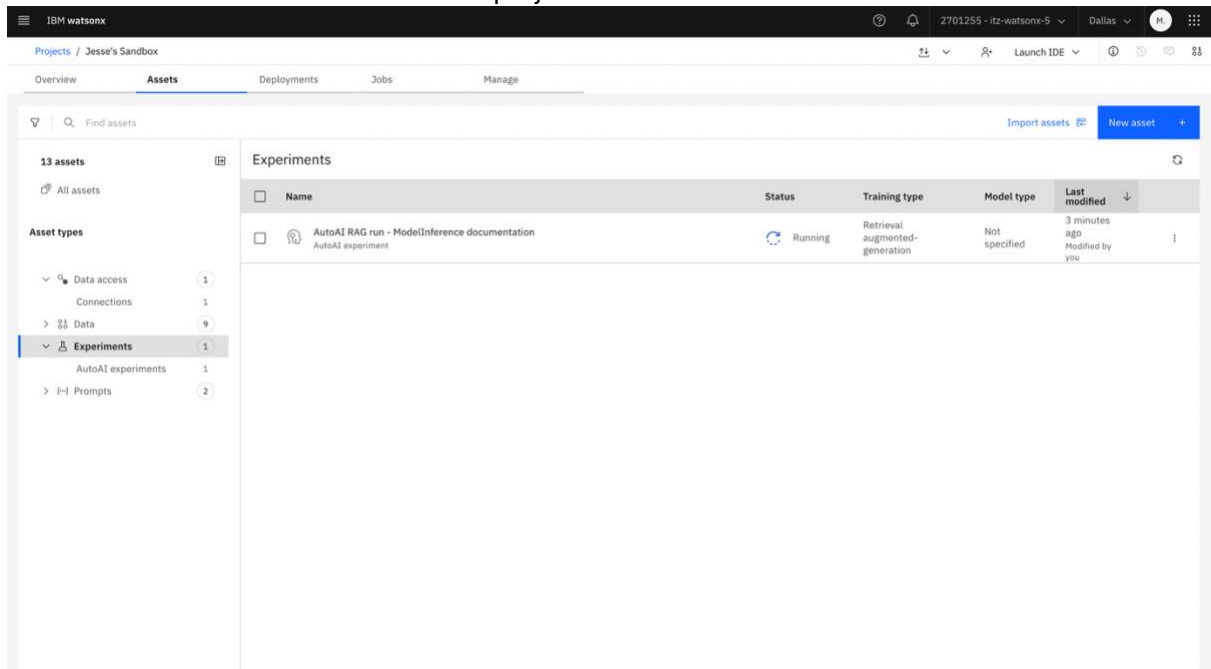10. In the IBM watsonx console sidebar, click Projects.

11. Click the New project button in the top-right.

12. Ensure the New tab is selected.

13. In Name, enter <Your Name>'s Sandbox.

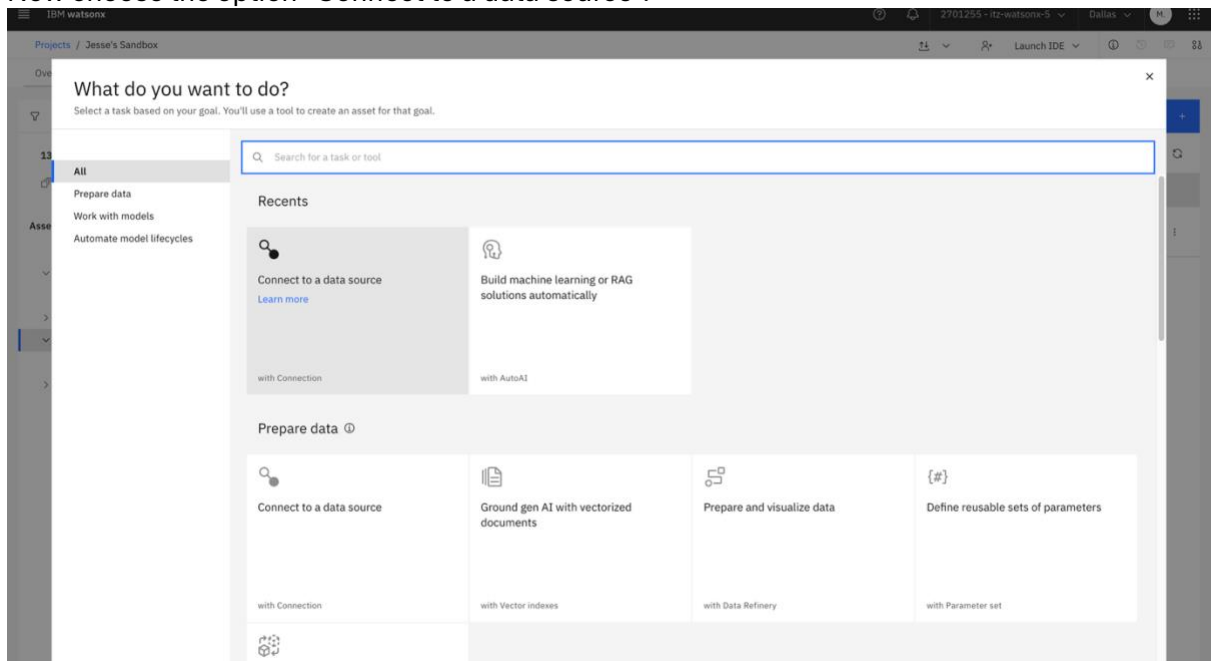14. (Optional) Fill in Description and Tags as desired.





Step 0.5: Import Elasticsearch Connection

1. Click on the "New Asset" button in the project.



2. Now choose the option "Connect to a data source".

3.  Select Elasticsearch from the options on left panel. Then click on Next.



4.  In a separate tab, navigate to the Databases for Elasticsearch service on IBM Cloud and copy the hostname, port, and TLS certificate under the Endpoints section of the Overview tab.



Copy the TLS certificate and add it to Elasticsearch connection on the watsonx.ai platform.

5. From the HTTPS tab you can note down the Hostname and Port number.
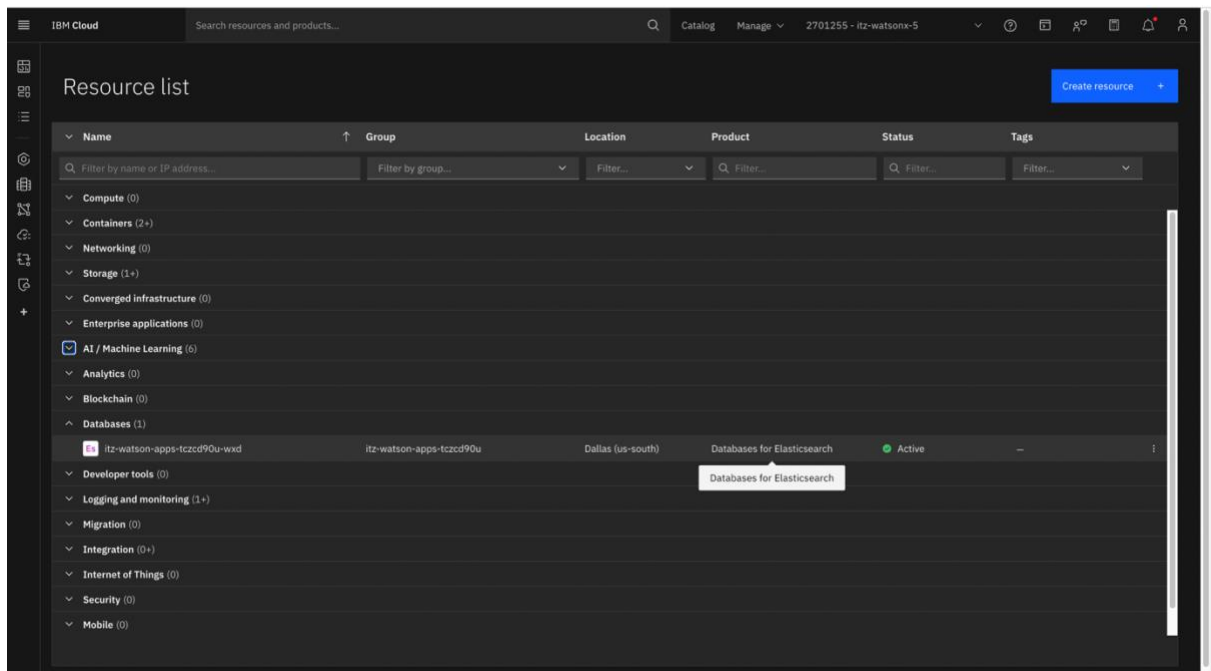


6. Go to the service credentials tab of the service, click on the arrow to expand it. Copy the username and password under "**connection.https.authentication**" in the service credential JSON.

7. Return to the project connection and set the following fields. Give the connection a name of: "Elasticsearch connection"



8. The URL should be formatted as https://{HOSTNAME}:{PORT}.

9. Now for authentication, from the dropdown select Username and password.



10. The username and password should be the ones copied from the service credentials.
11. The SSL certificate should be the TLS certificate.
12. Click on Test connection on the top, if that is successful then click on create.

Step 0.5: Import Data Asset

1. Navigate to the GitHub repo at

   https://github.com/JesseChan5171/JUCC-workshop#

2. Click Code → Download ZIP (or run git clone)
3. In the extracted directory, locate the "evaluation_dataset" folder.
4. In your project, click the **Assets** tab.
5. Click **Import assets** in the top-right.

6.  In the Import assets dialog, select **Local file** in the left menu.
7.  Click the **Data asset** tile.



8.  Click **Browse**, navigate to and select your dataset file (e.g. my_data.csv).
9.  Click **Import** to upload the all data file from "**evaluation_dataset**" folder (Step **3**).

**Step 1: Select the Task**

1. Navigate to the Assets Page inside projects in IBM watsonx.

2. Select **"Work with models"** from the list of categories.

3. Click on **"Build machine learning or RAG solutions automatically"** using the AutoAI tool.

**Step 2: Choose the Experiment Type**

1.  In the **AutoAI** experiment setup screen, choose what you want to build:

    o   Select **"Build a RAG solution"** to prepare and optimize RAG patterns based on document collections.

2.  Click **Next** to proceed.

**Step 3: Define Experiment Details**

1. Provide a name for your experiment (e.g., *AutoAI RAG run - ModelInference documentation*).

2. Optionally, add a description and tags for better organization and searchability.

3. In the configuration section:

   o Select the **watsonx.ai Runtime** service instance.

   o Choose the default environment definition

4. Click **Create** to proceed to the next step.

**Step 4: Add Document Sources**

1. In the **Add sources** screen, upload or select the required documents for the RAG solution:

    o Click **Select from project**, choose **Data asset** in the left-hand Categories pane to upload files from your project.

2. Select the upload all files and click **Open** to upload it.

Projects / Jesse_Sandbox / AutoAI RAG run - ModelInference documentation

## Select from project

Select which asset to use for your data collection and evaluation sources.

Search Categories

**Categories** 2

🔗 Connection  >
Data asset  >

Search Data assets

**Data assets** 20

- dataview.html
- eval_data.json
- flow_properties.html
- flow_scripting_example.html
- flow_scripting.html
- getting-started.html
- migration.html
- models-overview.html
- parameters.html
- reference_guides.html
- save_model.html
- spss_algorithms.html
- spss_tips.html
- spss_troubleshooting.html
- spss-comments.html
- spss-connections.html
- spss-modeler.html
- spss-modeling-process.html

### Select one or more Data assets

Select the asset you want to add.

Cancel    Select asset

**Step 5: Configure Indexing and Evaluation Sources**

1. In the **Configure details** page:

    o For indexing documents:

        ▪ Choose the **'Elasticsearch connection'** to store the vector index.

2. Once the source is uploaded, confirm its inclusion in the experiment.

**Step 6: Configure Indexing and Experiment Settings**

1.  Ensure that the documents and evaluation sources are listed correctly in the **Configure details** page.

2.  Click **Experiment settings**:

In **Experiment settings**:

- o Review the **Optimized metric** (**Answer faithfulness**) to prioritize the quality of generated answers.

- o Review the **Retrieval method** (**Window** or **Simple**) for data retrieval.

- o Review the **Foundation models** (llama-2-70b-instruct, flan-ul2-20b) for experimentation.

2. Save the changes.

**Step 7: Run the Experiment**

1. Return to the **Configure details** page.

2. Verify all configurations and sources, then click **Run experiment**.

3. Monitor the execution through the **Progress map**.

**Step 8: Review Results**

1. Navigate to the **Experiment summary**:

   - View the **Progress map** to confirm the experiment stages were completed successfully.

   - Check the **RAG Pattern leaderboard** for ranked results based on performance metrics such as **Answer faithfulness** and **Chunk size**.

2. Identify the top-performing pattern, such as Pattern 7 with llama-2-70b-instruct.

**Step 9: Pattern Comparison**

1.  Switch to the **Pattern comparison** tab for detailed insights:

    o   Use the **Metric chart** to compare patterns based on metrics like **Answer faithfulness** and **Correctness**.

    o   Analyze the trends to validate the optimal configuration for deployment.

**Step 10: Save the Results as a Notebook**

1.  After the experiment is completed:

    o   Review the **RAG Pattern leaderboard** to confirm the top-performing pattern.

2.  Click on the **Save as** button to generate a detailed notebook of the experiment results.

3. In the **Save notebook** dialog:

   - Under Choose your objective, select **Index building** (and Retrieval and generation if you want both).

   - Provide a name for the notebook (e.g., *AutoAI RAG - Top-performing Pattern 7*).

   - Choose the runtime environment and add an optional description or tags for easier identification.

4. Click **Create** to save the notebook for further use or sharing.

**Step 11: Open and Explore the Notebook**

1. Once the notebook is created:

    o Open the notebook to explore the detailed steps, data set, configurations, and metrics used in the experiment.



2. Run the notebook with following:

    o **Run the notebook cells**, check Instructions for setting up the vector store and index building.

    o **Watsonx.ai API key** for connection **(step 0)**

    o Details on RAG retrieval and generation processes.

    o Code snippets to reproduce the experiment and generate results.

3. Use the notebook to further refine or replicate the experiment as needed.

**Step 12: Extract RAG Pattern Configuration Keys**

1. Open the code or notebook cell containing your rag_pattern = { ... } definition.

```
The RAG Pattern details.

In [ ]: rag_pattern = {
            "composition_steps": [
                "model_selection",
                "chunking",
                "embeddings",
                "retrieval",
                "generation",
            ],
            "duration_seconds": 12,
            "name": "Pattern1",
            "settings": {
                "chunking": {"method": "recursive", "chunk_size": 1024, "chunk_overlap": 256},
                "embeddings": {
                    "truncate_strategy": "left",
                    "truncate_input_tokens": 512,
                    "model_id": "ibm/slate-125m-english-rtrvr",
                },
                "vector_store": {
                    "datasource_type": "elasticsearch",
                    "index_name": "autoai_rag_69210801_20250527023933",
                    "distance_metric": "cosine",
                    "operation": "upsert",
                    "schema": {
                        "id": "autoai_rag_1.1",
                        "name": "Document schema using open-source loaders",
                        "type": "struct",
                        "fields": [
                            {
                                "name": "pk",
                                "description": "Primary key",
                                "type": "string",
                                "role": "pk",
                            },
                            {
                                "name": "text_field",
                                "description": "text field",
                                "type": "string",
                                "role": "text",
                            },
                            {
                                "name": "document_id",
                                "description": "document name field",
                                "type": "string",
                                "role": "document_name",
                            },
                            {
                                "name": "start_index",
                                "description": "chunk starting token position in the source document",
                                "type": "number",
                                "role": "start_index",
```

2. Locate the top-level keys and the nested settings that drive each module.
3. Copy out the following key names for use in your module in later lab:

**IBM watsonx.ai Embeddings Key**

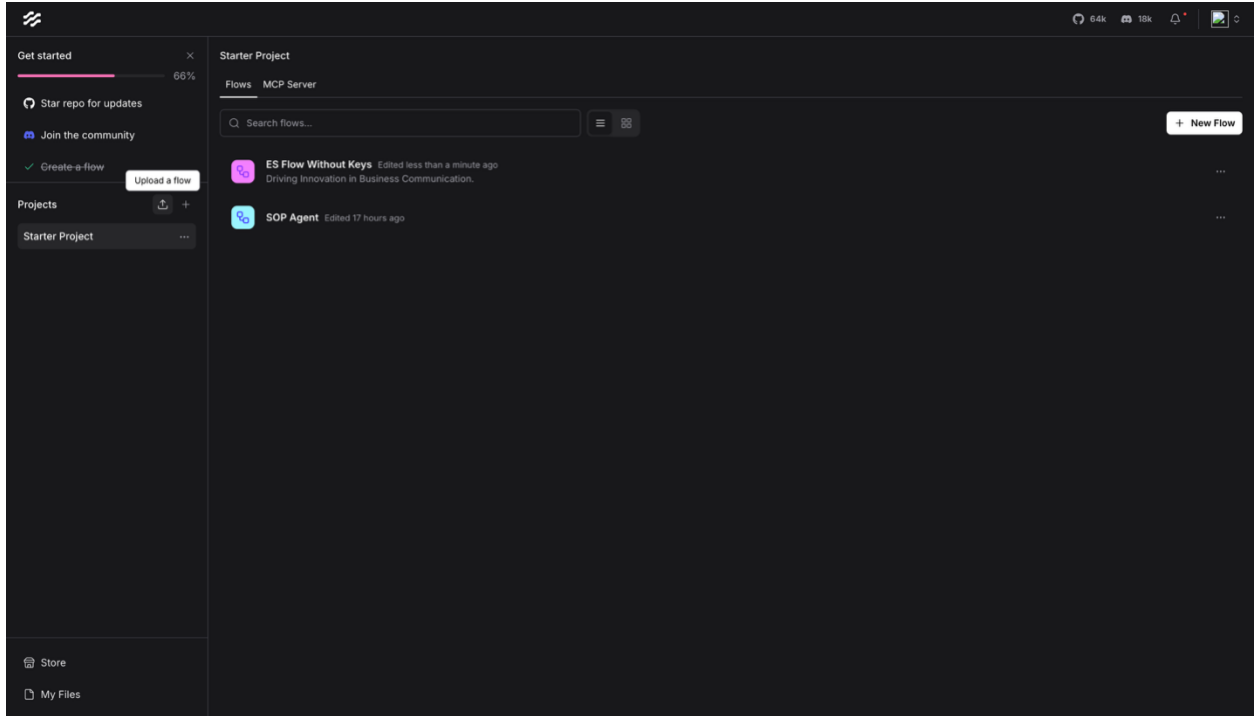| | Where to get it? | Value |
|---|---|---|
| watsonx API Endpoint | -> | https://us-south.ml.cloud.ibm.com |
| API Key | Step 0 | (get your own config) |
| Project ID | Find 'experiment_metadata' | (get your own config)<br><br>e.g. cd2dadd2-68e5-4680-8d79-cf8db08cc8c8 |
| Embedding Model Name | rag_pattern: embeddings -> model_id | (get your own config)<br><br>e.g. intfloat/multilingual-e5-large |

**IBM watsonx.ai Key**

| | Where to get it? | Value |
|---|---|---|
| watsonx API Endpoint | -> | https://us-south.ml.cloud.ibm.com |
| Model Name | rag_pattern: generation -> model_id | (get your own config)<br><br>e.g. meta-llama/llama-3-3-70b-instruct |

**Elasticsearch Vector Store**

| | Where to get it? | Value |
|---|---|---|
| Elasticsearch URL | -> | https://us-south.ml.cloud.ibm.com |
| Index Name | rag_pattern: vector_store -> index_name | (get your own config)<br><br>e.g. autoai_rag_69210801_20250527023933 |
| Username | Step 0 -> sub-step 6 | (get your own config)<br><br>e.g. ibm_cloud_c06543e4_08e4_42d9_836e_9de74a8cff8d |
| Password | Step 0 -> sub-step 6 | (get your own config)<br><br>e.g. eGq3GJdEh9kHhQHfetUu3gOE6hwGIP7l |

**Step 13: Upload the ES Flow into Lang flow**

1. In your browser, navigate to the Langflow UI: http://158.176.6.193:7860/
2. In the left-hand sidebar, click Upload a flow.
3. browse to the local JUCC-Workshop folder you downloaded earlier, select **ES Flow Without Keys.json**, and click Open.
4. Langflow will import the flow and display it on the canvas for you to review and configure
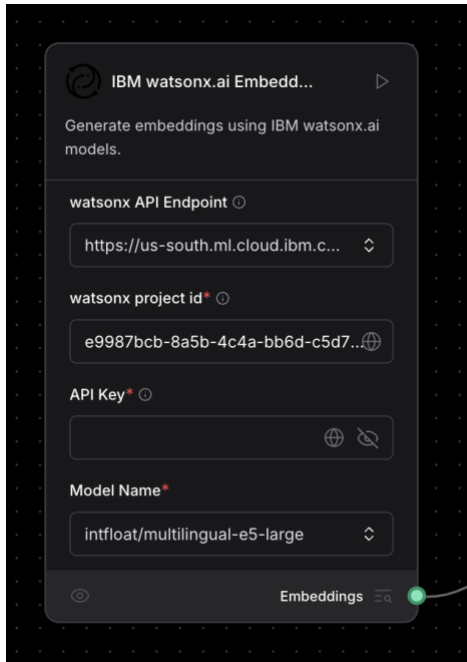
**Step 14: Configure Each Langflow Component and Understand Its Role**

1. **watsonx.ai Embeddings**
   Role: Transforms text (documents or queries) into vector representations that capture semantic meaning, enabling similarity-based search in later steps.
   - watsonx API Endpoint: https://us-south.ml.cloud.ibm.com
   - watsonx project id: your project ID
   - API Key: your API key
   - Model Name: your embeddings model (e.g. intfloat/multilingual-e5-large)

2.  Elasticsearch Vector Store
    Role: Stores and indexes those embeddings; given a query vector, it retrieves the top most similar document chunks to serve as context.
    - Elasticsearch URL: your Elasticsearch URL
    - Index Name: your index name (e.g. autoai_rag_69210801_20250527023933)
    - Username: your Elasticsearch username
    - Password: your Elasticsearch password

3. IBM watsonx.ai (Generation)
   Role: Consumes the retrieved chunks plus the user's query to generate a coherent, context-grounded answer using a foundation language model.
   - watsonx API Endpoint: https://us-south.ml.cloud.ibm.com
   - Model Name: your generation model (e.g. meta-llama/llama-3-3-70b-instruct)

**Step 15: Test Your RAG Flow in the Playground**

1. In the top-right corner of Langflow, click **Playground**.
2. In the "New chat" prompt that appears, type your test question (for example, "What is SPSS?").
3. Press **Enter** or click **Send**.
4. Observe how the flow:
   o Converts your question into an embedding,
   o Retrieves the most relevant document chunks from Elasticsearch,
   o Feeds those chunks + your question into the generation model,
   o Returns a grounded answer in the chat pane.
5. Review the **Session** sidebar to inspect earlier queries and responses or start a fresh chat by clicking the ➕ next to **Chat**.

Session May 27 at 06:37:36

**User**
Waht is SPSS?

**AI** meta-llama/llama-3-3-70b-instruct

SPSS (Statistical Package for the Social Sciences) is a software suite used for statistical analysis. It is widely used by researchers, data analysts, and students for data management, data analysis, and data visualization. SPSS provides a range of procedures for data manipulation, statistical analysis, and reporting. Unfortunately, I do not have any relevant documents to provide more information about SPSS. If you need more details, I suggest searching online or checking the official IBM SPSS website for more information.

Send