# watsonx Orchestrate Lab 2: Building a virtual assistant and setting up conversational search

## Use Case

In this lab you will learn how to set up the out-of-the-box conversational search using the Retrieval-Augmented Generation (RAG) pattern by connecting the assistant to an existing knowledge base (Elasticsearch) to search for relevant FAQs and generate an answer to the user's query using an LLM in watsonx.ai. You will also learn how to configure your assistant to provide general-purpose answering using an LLM and how to set up AI-guided actions to answer questions about a specific subject using provided knowledge.

The high-level steps to accomplish this are as follows:

1. Create a new assistant builder instance

2. Configure the conversational search extension in the assistant instance

3. Test the virtual assistant's content-grounded answering by asking questions contained in the FAQ document (stored in the Elasticsearch knowledge base) as well as additional questions answered by LLM only (general-purpose answering).

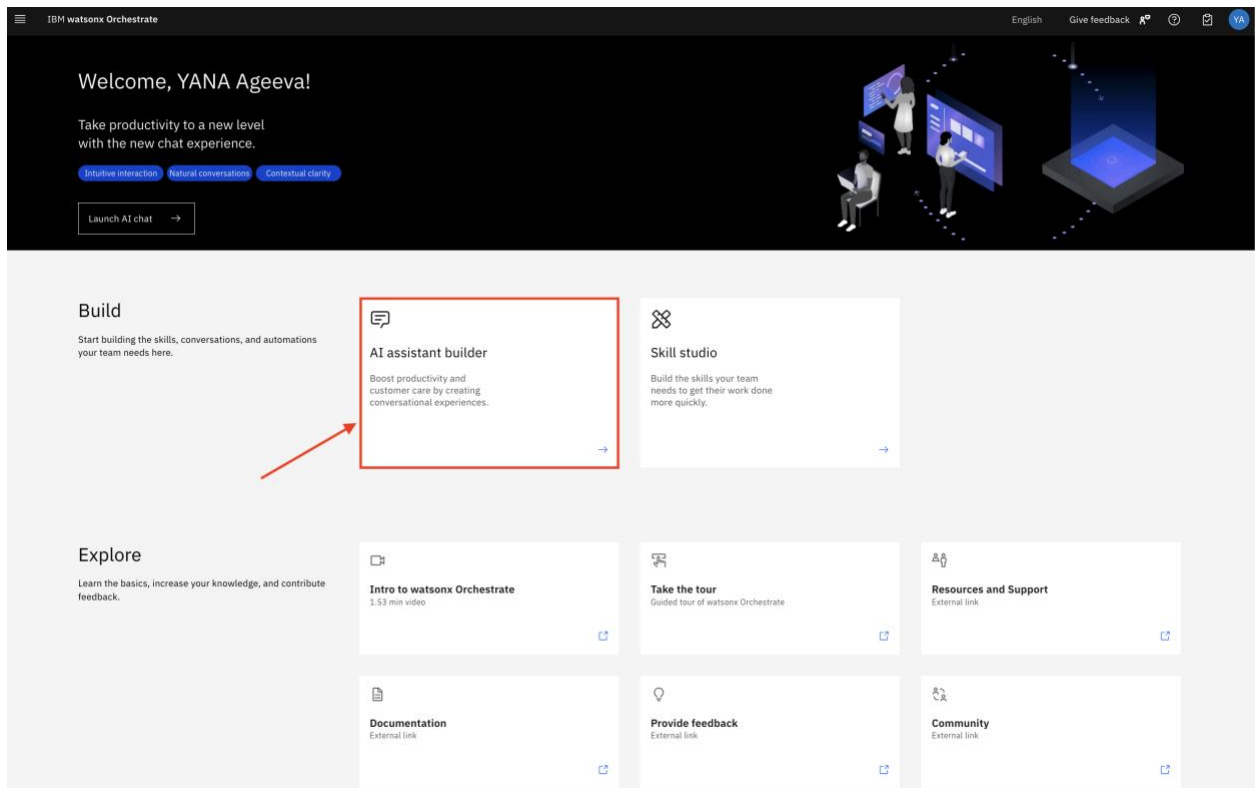4. Create and test an AI-guided action

**Assumptions**

1. You have access to an IBM watsonx Orchestrate tenant (on AWS or IBM Cloud)

2. To log in, use your IBM ID and password

3. You have access to an Elasticsearch instance

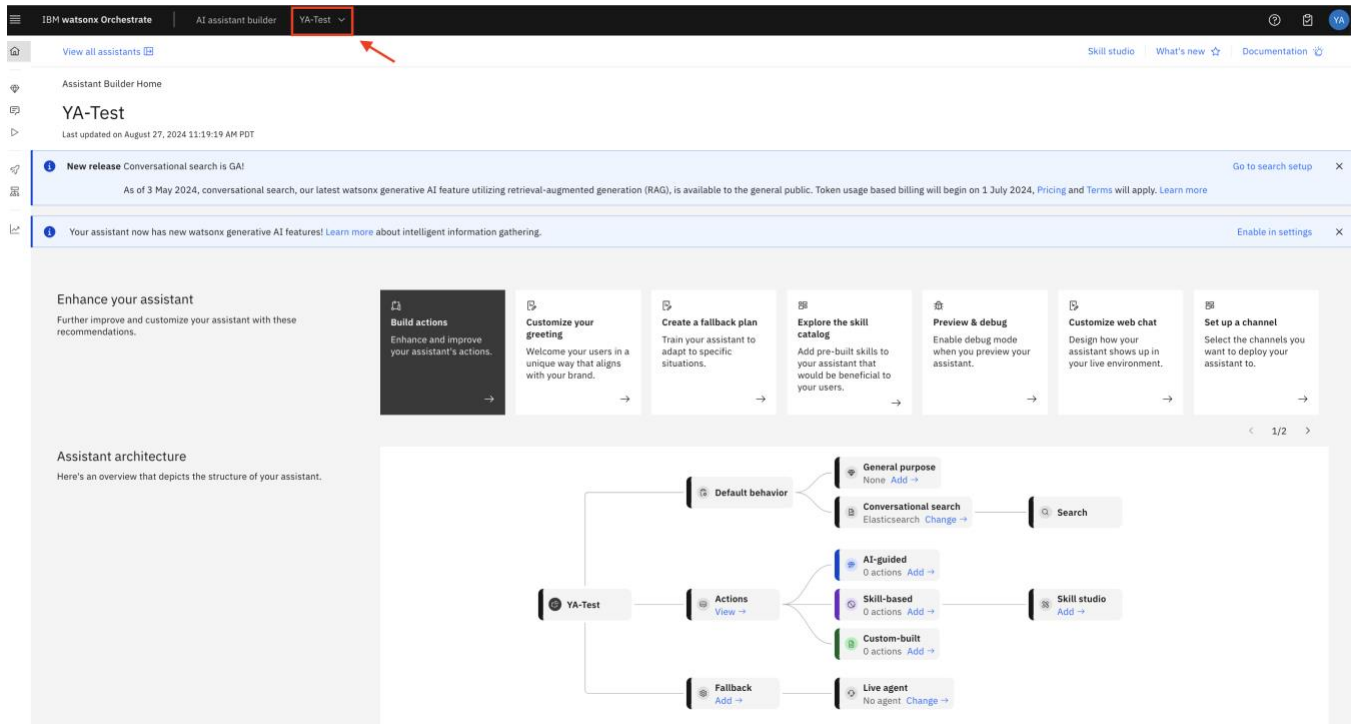4. Elasticsearch URL, port, username, and password

# Configure out-of-the-box conversational search using Elasticsearch

## Create a new virtual assistant

1. Using the environment that is provided to you, log in to watsonx Orchestrate.

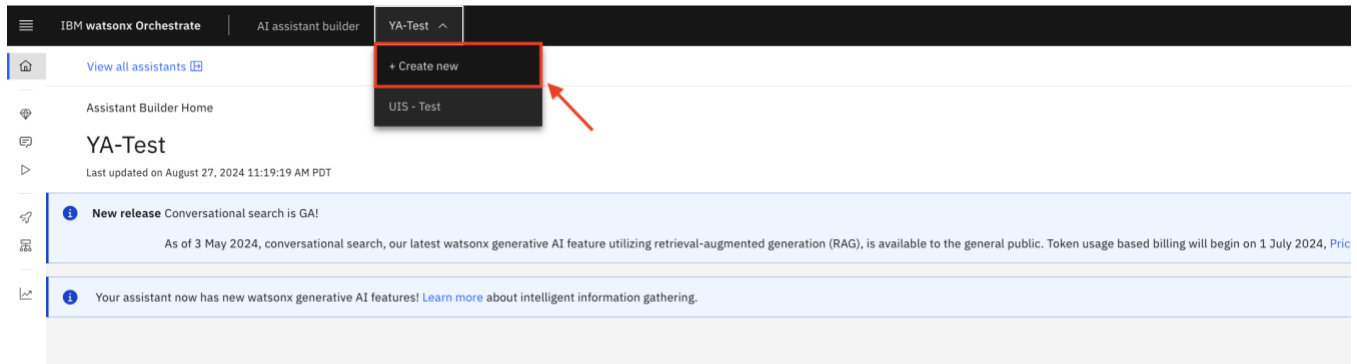2. Click on the Assistants tile to open the Assistant Builder:



3. You should see the following screen for one of the existing assistant instances. Click on the dropdown menu at the top to see existing instances:

**Note***: if there are no existing instances, you will have to create the first one (answer the questions in the wizard and use default settings when possible.*

4. Create a new assistant instance



5. Name your assistant as FirstLast - Assistant (where First is your first name, Last is your last name) and click **Create assistant**.

Create a new assistant

Assistant name

FirstLast

Your assistant name will be kept internally and not visible to your customers

Description (optional)                                                                                    0/128

Add a description for this assistant

Assistant language

English (US)

This is the language your assistant will speak.

Cancel                                                                        Create assistant

6.  Your assistant has been created and you can now see it in the dropdown list.



Congratulations! You have successfully created your first assistant. It is connected by default to the IBM Granite LLM for general purpose answering, and you can test it by asking questions in the preview chat window.
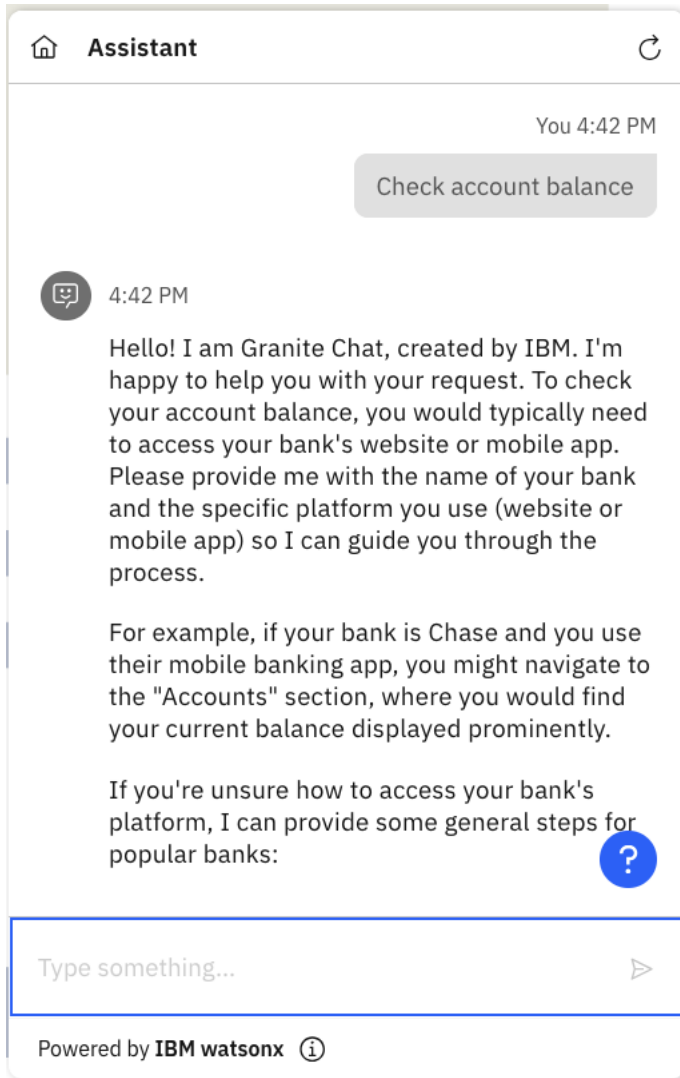
7. Click on **Preview** and move your cursor into the text area of the assistant that pops up:



8. Enter any query e.g. "Check account balance" (or you can also click on one of the default buttons in the chat window)

**Assistant**

## Hi! I'm a virtual assistant. How can I help you today?

Example: Find nearby location

Example: Check account balance ←

Example: See how I can help

**1** **2**

Check account balance ➤
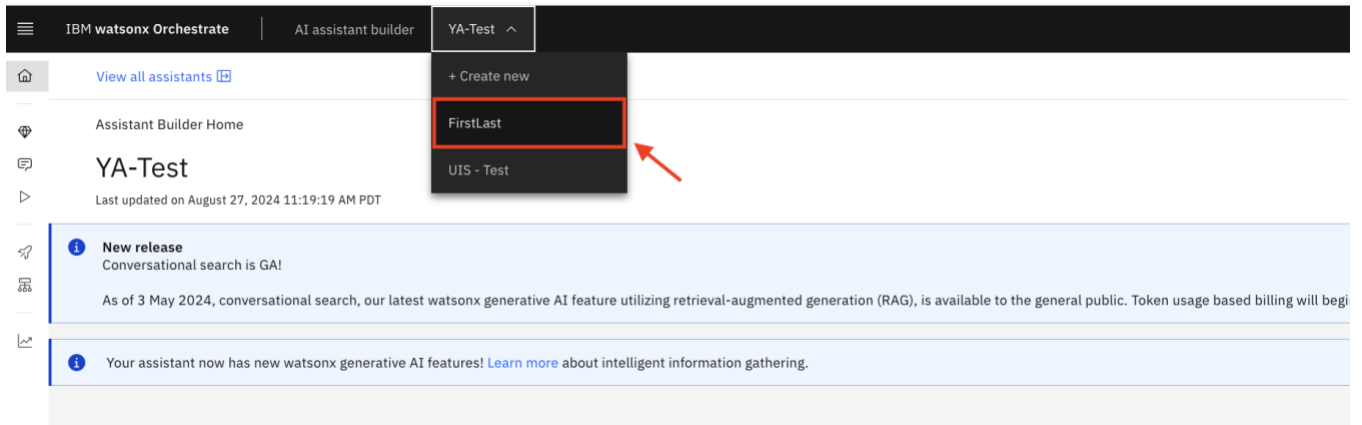
Powered by **IBM watsonx** ⓘ

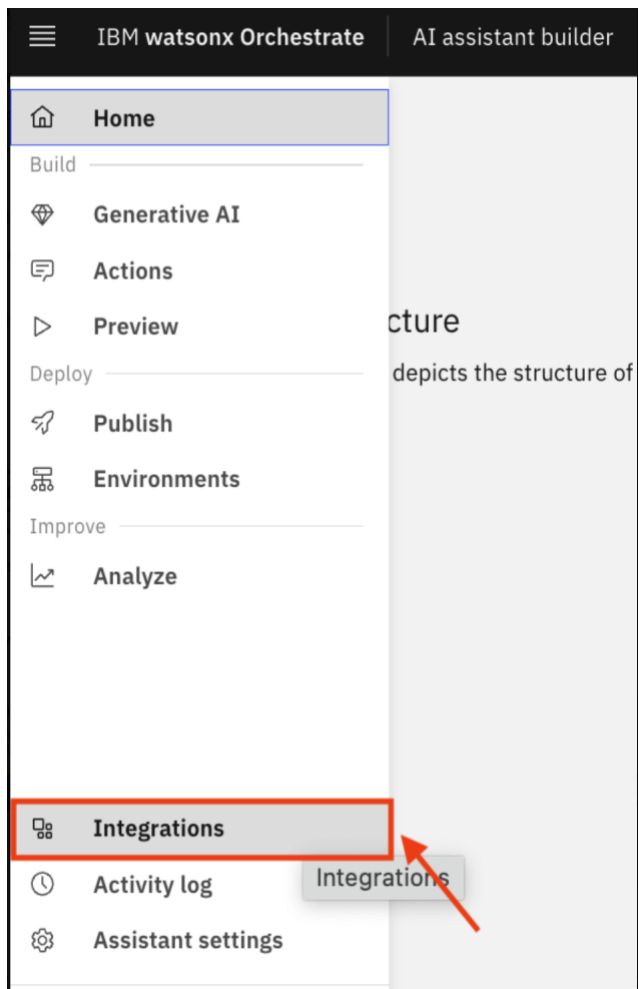9. The answer is provided by the Granite LLM using general purpose answering:

As you can see, this is a very generic answer because the assistant doesn't know anything about you or which bank you have an account at. In the following steps we will set up the assistant and teach it to answer questions intelligently, based on the provided knowledge base and implemented actions.

# Set up conversational search in your assistant builder instance

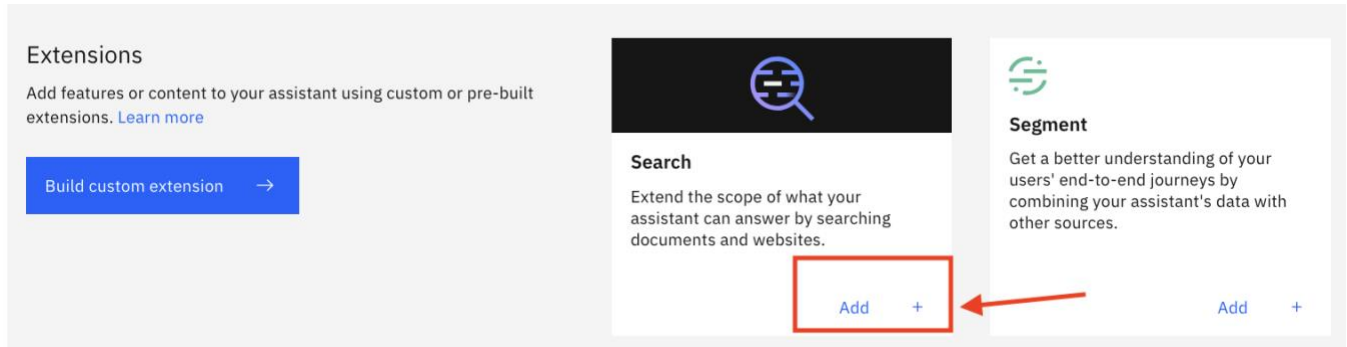1. Select your assistant instance if it isn't selected already
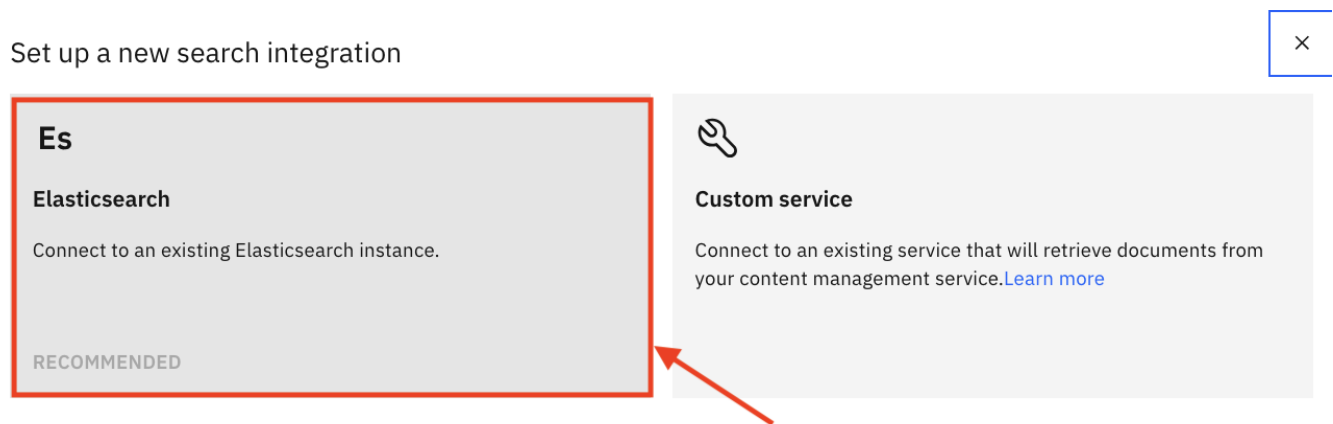
2. Open the **Integrations** tab

3.  Scroll down to **Extensions** and add **Search**



4.  Select **Elasticsearch**



**Note**: you can also connect to watsonx.data Milvus instead of Elasticsearch – follow this guide.

5.  You should now see the following screen prompting you to enter **Elastcisearch** credentials:

## Elasticsearch  Draft

Close | Next

**○ Connect Elasticsearch**          ○ Select index          ○ Conversational search (optional)

## Connect to your Elasticsearch instance.

Fill in the information below to access your Elasticsearch instance.  Learn more

Elasticsearch url

```
exampleUrl.com
```

Elasticsearch port (optional)

```
123
```

Choose an authentication type

```
Basic authentication                          ⌄
```

Elasticsearch username

```
exampleUser
```

Elasticsearch password

```
••••••••••                                    👁
```

6. Connect to the Elasticsearch instance in watsonx Orchestrate using the following values provided by your proctor:

   **Elasticsearch url** (make sure to remove the trailing slash, if any)

   **Elasticsearch port**

Leave authentication type as **Basic authentication** and fill in the **Elasticsearch username** and **Elasticsearch password**



7. Click **Next**

8. On the next screen select **Upload documents to a new index in your Elasticsearch instance** to automatically create a new index in Elasticsearch and click **Next** to continue

9. On the next screen make sure conversational search is **On** and click **Save**



Once you click **Save,** you will be taken to the **Documents** tab where you can upload your documents to the knowledge base (see next section).

# Upload documents to the knowledge base

1. Now we need to upload our documents into the knowledge base. In the **Documents** tab click **Upload**



2. Find the file  **New Windows 11 Known Issues Keep Coming.pdf** (provided by your proctor) and drag and drop the .pdf file into the drag and drop area. Click on **Upload files**



3. Once the document is ready, the screen will look the following way and you can click on **Close**

And this is it. You have now finished the process of ingesting documents into the **Elasticsearch** knowledge base!

# Test the virtual assistant

Now that conversational search is set up and the document has been uploaded to the knowledge base, we can test the virtual assistant by asking some questions 1) based on the Q&As in the document 2) general questions that may be answered by a built in LLM.

1.  Open **Preview** and go into the input window in the virtual assistant:

2. Type a question that is answered in the document, e.g. "how do I roll back from windows 11 to windows 10?" and check that the answer provided was retrieved from the document:

## Assistant

6:36 AM

To roll back from Windows 11 to Windows 10, follow these steps:

1. Go to Settings, click System and Recovery.
2. Click Go back next to the Previous version of Windows.
3. Check the reasons why you want to uninstall when prompted.
4. Click No, thanks when asked to check for updates instead of rolling back.
5. When a "What you need to know" window pops up, click Next to continue.
6. Click Next again when it warns you to remember your password from Windows 10.
7. Click Go back to an earlier build.
8. The system will now restart and go through a restoration process that should take a few minutes. When that completes, you're back on Windows 10.

Type something...

Powered by **IBM watsonx** ⓘ

3. Type a question that is NOT answered in the document, e.g. "Are there lions in Yosemite?" and check if the LLM provides an answer

AI  Accuracy of generated answers may vary.

You 6:38 AM

Are there lions in Yosemite?

6:38 AM

While Yosemite National Park is home to a variety of wildlife, including black bears, elk, and coyotes, it does not have a lion population. Lions are native to Africa and the Middle East, and there are no known wild lions living within the boundaries of Yosemite National Park or any other national parks in North America.

I hope this information is helpful! Let me know if you have any other questions.

Type something...                                       ▷

Powered by **IBM watsonx** ⓘ

4. **(Optional)** You can test out different LLMs available under the **Generative AI** tab and observe how different the answers are to the same question

5. **(Optional)** You can also experiment with the answer behavior by toggling **General-purpose answering** and **Content-grounded answering** options on and off. See how the answers change with different combinations of these settings.

**Note:** Given the timeframe available for this hands-on lab, we experimented with only one .pdf file. However, you could upload multiple files and ask questions about information contained in any of them. The answers provided by conversational search will be reference the original source.
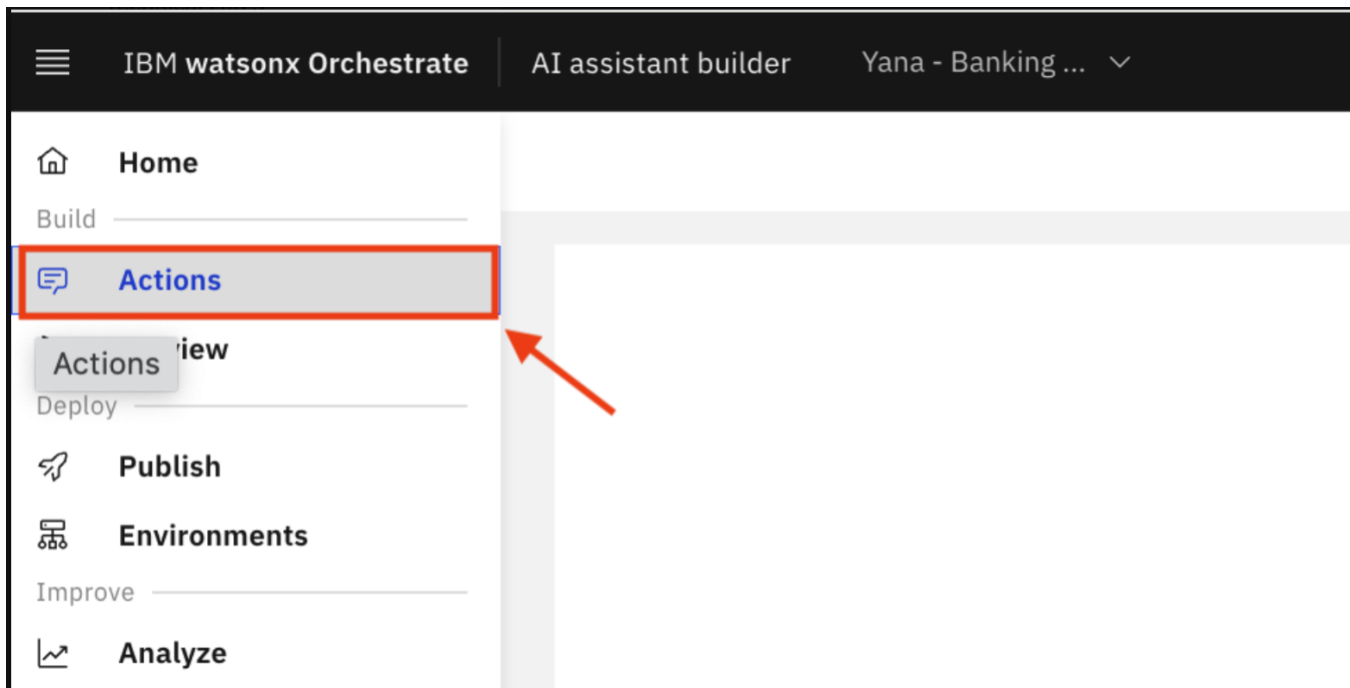
In the steps above you configured conversational search in your assistant with a knowledge base that stores FAQs (in Elasticsearch). When the user asks a question, the knowledge base is queried to retrieve any relevant FAQs which are then passed to a built-in watsonx.ai LLM (IBM Granite) to generate an answer for the customer.

# Implement an AI guided action

## Create an AI-guided action to answer questions

In this step we will create an **AI-guided action** to answer questions from a user related to a specific subject, given provided knowledge context. This is different from conversational search in that the answers are grounded in very specific knowledge. At the same time, the answers to AI-guided actions are limited to the specific knowledge provided vs. the **general-purpose answering** where LLMs are allowed to answer the user's question based on all the vast data they have been trained on.

6.  Go into your assistant builder instance (make sure your name is selected in the drop-down) and select the **Actions** tab



7.  Create a new action and select **AI-guided action**

8. Pick a topic that the user may want to ask about, e.g. IT Service Management (ITSM) and provide a phrase that a user might type that is related to this subject, e.g. "What is ITSM"



9. Go into **Customer starts with:** in the action configuration menu and provide additional phrases/words that the user might type to ask about this subject, each time hitting Enter, and then **Save** everything

10. Identify related content on this subject and provide it as **knowledge** inside the action (you can use the text below the following screenshot as a sample). Make sure to save again and then click on the preview button in the bottom right of the screen to test your action



## Text that can be used as sample knowledge:

**Understanding IT Service Management (ITSM)**

**IT Service Management (ITSM)** refers to the entirety of activities, policies, processes, and procedures that organizations use to design, deliver, manage, and improve the way IT services are used within their business. ITSM ensures that the right mix of people, processes, and technology are in place to provide value to customers.

**Key Components of ITSM**

1. **Incident Management**: Restoring normal service operation as quickly as possible.

2. **Problem Management**: Identifying the root cause of incidents and finding permanent solutions.

3. **Change Management**: Managing changes to IT services with minimal disruption.

4. **Service Request Management**: Handling requests from users for information, advice, or access to IT services.

5. **Asset and Configuration Management**: Keeping track of IT assets and configurations.

6. **Knowledge Management**: Creating, sharing, using, and managing knowledge and information.

**ITSM Workflows for Creating a Ticket**

Creating a ticket in ITSM involves several steps. Here's a generalized workflow:

1. **Ticket Creation**

   - **User Submission**: A user submits a ticket via an online portal, email, or phone.

   - **Automatic Creation**: The system may create tickets automatically through monitoring tools.

2. **Categorization**

   - **Ticket Classification**: Categorize the ticket based on the type of issue (e.g., incident, service request).

   - **Priority Assignment**: Assign priority based on the impact and urgency of the issue.

3. **Assignment**

   - **Routing**: Route the ticket to the appropriate support group or technician.

   - **Automatic Assignment**: The system may automatically assign tickets based on predefined rules.

4. **Investigation and Diagnosis**

   - **Initial Assessment**: The assigned technician reviews the ticket details.

   - **Diagnosis**: Investigate the issue and diagnose the root cause.

   - **Update Ticket**: Update the ticket with diagnostic information.

5. **Resolution and Recovery**

   - **Fix Implementation**: Implement the solution or workaround.

   - **User Communication**: Communicate the solution to the user and confirm resolution.

   - **Update Ticket**: Document the steps taken to resolve the issue.

6. **Closure**

   - **User Confirmation**: Confirm with the user that the issue is resolved satisfactorily.

   - **Ticket Closure**: Close the ticket and update the system.

7. **Post-Closure Review**

   - **Review and Analysis**: Conduct a review to analyze the issue and identify improvements.

   - **Knowledge Base Update**: Update the knowledge base with new insights or solutions.

**Example of an ITSM Ticket Workflow**

1. **Step 1**: Ticket Creation

   - **User**: John Doe reports an issue with email access through the IT service portal.

2. **Step 2**: Categorization

   - **System**: Categorizes the ticket as an "Email Issue" with a high priority.

3. **Step 3**: Assignment

   - **System**: Automatically assigns the ticket to the "Email Support" team based on predefined rules.

4. **Step 4**: Investigation and Diagnosis

   - **Technician**: Jane Smith reviews the ticket and identifies a server issue as the root cause.

- **Technician**: Documents her findings and the steps she will take to resolve the issue.

5. **Step 5**: Resolution and Recovery

    - **Technician**: Implements the fix by restarting the email server.

    - **Technician**: Notifies John Doe that the issue has been resolved and asks him to verify.

6. **Step 6**: Closure

    - **User**: John Doe confirms that the email access issue is resolved.

    - **Technician**: Closes the ticket and updates the resolution details in the system.

7. **Step 7**: Post-Closure Review

    - **Team Lead**: Conducts a review of the incident to identify any improvements or patterns.

8. **Knowledge Base**: Updates the knowledge base with the resolution steps for future reference.

This workflow ensures efficient handling of issues, leading to better service quality and user satisfaction.

# Test your AI guided action

1. Test the action by typing some of the sample phrases/words you provided earlier and make sure your new action was invoked and the results are based on the knowledge you provided:

**Preview**  ↻  —

how do I create a ticket?

4:35 PM

Creating a ticket in ITSM involves several steps. Here's a generalized workflow:

1. Ticket Creation
   User Submission: A user submits a ticket via an online portal, email, or phone.
   Automatic Creation: The system may create tickets automatically through monitoring tools.
2. Categorization
   Ticket Classification: Categorize the ticket based on the type of issue (e.g., incident, service request).
   Priority Assignment: Assign priority based on the impact and urgency of the issue.
3. Assignment
   Routing: Route the ticket to the appropriate support group or technician.
   Automatic Assignment: The system may automatically assign tickets based on predefined rules.
4. Investigation and Diagnosis
   Initial Assessment: The assigned technician reviews the ticket details.

Type something...  ▷

2. Feel free to experiment with additional settings, e.g. selecting a different LLM

This final step concludes the lab. You just configured and tested conversational search to provide content-grounded answers to users' questions. You also experimented with general-purpose answering using LLMs. Finally, you created an AI-guided action to answer questions from a user related to a specific subject, given provided knowledge.