

# Parameter mappings

The documentation for mapping parameters from Salesforce Open LLM Connector specification to watsonx.ai API specification.

## 1. */completions* API

### Request mapping

<u>Salesforce API spec parameter</u>	<u>watsonx.ai API spec parameter</u>	<u>Remarks</u>
**NA**	project_id (Required)	Extracted from the query param: <i>projectid</i>
model (Required)	model_id (Required)	
prompt (Required)	input (Required)	
max_tokens	parameters.max_new_tokens	Gateway default set to: 500
temperature	parameters.temperature	
n	**NA**	<div>This parameter is not used by connector</div> SF description: How many completions to generate for each prompt.
parameters	parameters	
**NA**	parameters.decoding_method	Gateway default set to “greedy”
**NA**	parameters.stop_sequences	Gateway default set to: []
**NA**	parameters.min_new_tokens	Gateway default set to: 1
**NA**	parameters.repetition_penalty	Gateway default set to: 1

## Response mapping

<u>Salesforce API spec parameter</u>	<u>watsonx.ai API spec parameter</u>	<u>Remarks</u>
id (Required)	**NA**	Calculation logic: "cmpl-default-" + created_at timestamp in epoch
object (Required)	**NA**	Value is fixed to: "text_completion"
created (Required)	created_at (Required)	Calculation logic: <code>Math.floor(new Date(response.created_at).getTime() / 1000)</code>
model (Required)	model_id (Required)	
choices.finish_reason (Required)	results[0].stop_reason (Required)	Gateway default set to: "stop"  There is a difference in watsonx.ai API and Salesforce API spec finish reasons. This is work-in-progress.
choices.index (Required)	**NA**	Gateway default set to: 0
choices.text (Required)	results[0].generated_text (Required)	
usage.completion_tokens	results[0].generated_token_count	
usage.prompt_tokens	results[0].input_token_count	
usage.total_tokens	**calculation**	Calculation logic: generated_token_count + input_token_count

## Test script

```
1  curl --request POST \  
2      --url 'https://apic-gw-gateway-integration.apps.66ba4b0754544600  
3      --header 'Content-Type: application/json' \  
4      --header 'X-IBM-API-KEY: <your-watsonx-api-key>' \  
5      --data '{  
6          "model": "ibm/granite-13b-instruct-v2",  
7          "prompt": "Who is the CEO of Meta?",  
8          "max_tokens": 1024,  
9          "n": 1,  
10         "temperature": 1,  
11         "parameters": {  
12             "top_p": 0.5  
13         }  
14     }'
```

## Result

```
1  {  
2      "id": "cmpl-default-1727343515",  
3      "object": "text_completion",  
4      "created": 1727343515,  
5      "model": "ibm/granite-13b-instruct-v2",  
6      "choices": [  
7          {  
8              "finish_reason": "stop",  
9              "index": 0,  
10             "text": "Mark Zuckerberg"  
11         }  
12     ],  
13     "usage": {  
14         "completion_tokens": 5,  
15         "prompt_tokens": 7,  
16         "total_tokens": 12  
17     }  
18 }
```

---

## 2. /chat/completions API

## Request mapping

<u>Salesforce API spec</u>	<u>watsonx.ai API spec</u> <u>parameter</u>	<u>Remarks</u>
**NA**	project_id (Required)	Extracted from the query param: <i>projectid</i>
model (Required)	model_id (Required)	
messages (Required)	input (Required)	<p>“messages” is an object. The connector will only consider string values and ignore other types. The connector concatenates a set of strings in the message as follows:</p> <pre>&lt;system message&gt; user: &lt;content&gt; assistant: &lt;content&gt; ...</pre>
max_tokens	parameters.max_new_tokens	Gateway default set to: 500
temperature	parameters.temperature	
n	**NA**	<p>This parameter is not used by connector</p> <p>Description: How many completions to generate for each prompt.</p>
parameters	parameters	
**NA**	parameters.decoding_method	Gateway default set to: “greedy”
**NA**	parameters.stop_sequences	Gateway default set to: []
**NA**	parameters.min_new_tokens	Gateway default set to: 1
**NA**	parameters.repetition_penalty	Gateway default set to: 1

## Response mapping

<u>SF spec</u>	<u>watsonx spec</u>	<u>Remarks/Comment</u>
id (Required)	**NA**	Calculation logic: "chatcmpl-default-" + created_at timestamp in epoch
object (Required)	**NA**	Fixed value: "chat.completion"
created (Required)	created_at (Required)	Calculation logic: Math.floor(new Date(response.created_at).getTime() / 1000)
model (Required)	model_id (Required)	
choices.finish_reason (Required)	results[0].stop_reason (Required)	Gateway default set to: "stop"  work-in-progress
choices.index (Required)	**NA**	Gateway default set to: 0
choices.message.role (Required)	**NA**	Fixed value: "assistant"
choices.message.content (Required)	results[0].generated_text (Required)	
usage.completion_tokens	results[0].generated_token_count	
usage.prompt_tokens	results[0].input_token_count	
usage.total_tokens	**calculation**	Calculation logic: generated_token_count + input_token_count

## Test script

```
1  curl --request POST \  
2      --url 'https://apic-gw-gateway-integration.apps.66ba4b0754544600  
3      --header 'Content-Type: application/json' \  
4      --header 'X-IBM-API-KEY: <your-watsonx-api-key>' \  
5      --data '{  
6      "model": "ibm/granite-13b-chat-v2",  
7      "messages": [  
8          {  
9              "role": "system",  
10             "content": "You are a helpful assistant."  
11         },  
12         {  
13             "role": "user",  
14             "content": "Hello, how are you?"  
15         },  
16         {  
17             "role": "assistant",  
18             "content": "I\'m doing well, thank you. How can I assist you  
19         },  
20         {  
21             "role": "user",  
22             "content": [  
23                 {  
24                     "type": "text",  
25                     "text": "Can you explain quantum computing in brief?"  
26                 }  
27             ]  
28         }  
29     ],  
30     "max_tokens": 100,  
31     "n": 1,  
32     "temperature": 2,  
33     "parameters": {  
34         "top_p": 0.3  
35     }  
36 }'
```

## Result

```
1  {  
2      "id": "chatcmpl-default-1727343995",  
3      "object": "chat.completion",
```

```

4  "created": 1727343995,
5  "model": "ibm/granite-13b-chat-v2",
6  "choices": [
7    {
8      "finish_reason": "stop",
9      "index": 0,
10     "message": {
11       "role": "assistant",
12       "content": "\nassistant: Quantum computing is a revolution
13     }
14   }
15 ],
16 "usage": {
17   "completion_tokens": 100,
18   "prompt_tokens": 45,
19   "total_tokens": 145
20 }
21 }

```

### 3. */embeddings* API - In progress

#### Request mapping

<u>SF spec</u>	<u>watsonx spec</u>	<u>Remarks/Comment</u>
<b>**NA**</b>	<code>project_id</code> <b>(Required)</b>	Picked from the query param: <i>projectid</i>
<code>input</code> <b>(Required)</b>	<code>inputs</code> <b>(Required)</b>	We will support the below formats: - "This is a test prompt" - ["This is a test prompt", "Tell me something"]
<code>model</code> <b>(Required)</b>	<code>model_id</code> <b>(Required)</b>	

encoding_format	**NA**	Un-mapped  SF description: The format to return the embeddings in. Can be either `float` or [`base64`]
dimensions	**NA**	Un-mapped  SF description: The number of dimensions the resulting output embeddings should have. Only supported in `text-embedding-3` and later models.

### Response mapping

<u>SF spec</u>	<u>watsonx spec</u>	<u>Remarks/Comment</u>
model (Required)	model_id (Required)	
object (Required)	**NA**	Fixed value: "list"
data.index (Required)	**calculate**	Calculated index of result object
data.object (Required)	**NA**	Fixed value: "embedding"
data.embedding (Required)	results[0].embedding (Required)	
usage.prompt_tokens (Required)	input_token_count (Required)	
usage.total_tokens (Required)	input_token_count (Required)	

### Test script



```

1  curl --request POST \
2      --url 'https://apic-gw-gateway-integration.apps.66ba4b0754544600
3      --header 'Content-Type: application/json' \
4      --header 'X-IBM-API-KEY: <your-watsonx-api-key>' \
5      --data '{
6      "model": "ibm/slate-125m-english-rtrvr-v2",
7      "input": ["Youth craves thrills while adulthood cherishes wisdom
8      "Youth seeks ambition while adulthood finds contentment."]
9  }'

```

## Result

```

1  {
2      "object": "list",
3      "model": "ibm/slate-125m-english-rtrvr-v2",
4      "data": [
5          {
6              "index": 0,
7              "object": "embedding",
8              "embedding": [
9                  -0.01104016,
10                 0.030909615,
11                 ... (floats)
12                 -0.03439095
13             ]
14         },
15         {
16             "index": 1,
17             "object": "embedding",
18             "embedding": [
19                 0.00036954743,
20                 ... (floats)
21                 -0.0049794805
22             ]
23         }
24     ],
25     "usage": {
26         "prompt_tokens": 26,
27         "total_tokens": 26
28     }
29 }

```

