

LLaMa

Example Generation

Future of Programming

WiSe 2023 / 2024

Software Architecture Group

Christian Raue / Lara Kursawe

**Design IT.
Create Knowledge.**

www.hpi.de



Question: Can LLMs generate helpful examples of API interfaces?



You

Show me an example of how to use the `DataFrame.empty` method from the pandas API

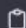


ChatGPT

The `DataFrame.empty` method in the pandas library is used to check if a DataFrame is empty (i.e., it has zero rows and zero columns). It returns a boolean value indicating whether the DataFrame is empty or not.

Here's an example of how to use the `DataFrame.empty` method:

python

 Copy code

```
import pandas as pd

# Create an empty DataFrame
empty_df = pd.DataFrame()

# Check if the DataFrame is empty
if empty_df.empty:
    print("The DataFrame is empty.")
else:
    print("The DataFrame is not empty.")
```

In this example, we first create an empty DataFrame using the `pd.DataFrame()` constructor. Then, we use the `empty` method to check if the DataFrame is empty. If it is empty, it will print "The DataFrame is empty," otherwise, it will print "The DataFrame is not empty."



You

Show me an example of how to use the `PreTrainedTokenizer.convert_ids_to_tokens` method from HuggingFace

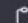


ChatGPT

As of my last knowledge update in January 2022, Hugging Face's Transformers library provides the `PreTrainedTokenizer.convert_ids_to_tokens` method for converting token IDs back to tokens. Note that the library may have evolved, and it's always a good idea to check the latest documentation for any changes or updates.

Here's an example of how you might use `convert_ids_to_tokens`:

python

 Copy code

```
from transformers import PreTrainedTokenizer, AutoTokenizer

# Replace 'bert-base-uncased' with the desired pre-trained model name
model_name = 'bert-base-uncased'
tokenizer = AutoTokenizer.from_pretrained(model_name)

# Encode a sentence to obtain token IDs
text = "Hello, how are you doing?"
token_ids = tokenizer.encode(text, add_special_tokens=True)

# Convert token IDs back to tokens
tokens = tokenizer.convert_ids_to_tokens(token_ids)

print("Original text:", text)
print("Token IDs:", token_ids)
print("Tokens:", tokens)
```



Question: Can LLMs generate helpful examples of API interfaces?

Yes!

Thank you for your attention!

[f convert_ids_to_tokens](#)[<source>](#)

```
( ids: typing.Union[int, typing.List[int]], skip_special_tokens: bool = False ) → str or List[str]
```

Parameters

- **ids** (int or List[int]) — The token id (or token ids) to convert to tokens.
- **skip_special_tokens** (bool, *optional*, defaults to False) — Whether or not to remove special tokens in the decoding.

Returns str or List[str]

The decoded token(s).

Converts a single index or a sequence of indices in a token or a sequence of tokens, using the vocabulary and added tokens.

[f convert_tokens_to_ids](#)[<source>](#)

```
( tokens: typing.Union[str, typing.List[str]] ) → int or List[int]
```

Parameters

- **tokens** (str or List[str]) — One or several token(s) to convert to token id(s).

Returns int or List[int]

The token id or list of token ids.

Converts a token string (or a sequence of tokens) in a single integer id (or a sequence of ids), using the vocabulary.

Many APIs don't provide any examples in their documentation, i.e. HuggingFace. This makes it difficult for inexperienced users to get familiar with the API, having to use other resources to help them find examples.

Idea

When using a less known API, you might even have to provide the API documentation for every function you need an example for.

Chatbots

- Helpful when working and learning new tools, but often cumbersome to use
 - Might need long prompts
 - First results may be unhelpful
 - Multiple follow up prompts can be annoying
- Different LLMs give different results (of different quality)

LA You

Show me an example of how to use the `PreTrainedTokenizer.convert_ids_to_tokens` method from HuggingFace



ChatGPT

As of my last knowledge update in January 2022, Hugging Face's Transformers library provides the `PreTrainedTokenizer.convert_ids_to_tokens` method for



You

Okay how would i use `get_added_vocab`?



ChatGPT

The `get_added_vocab` method in Huggingface's Transformers library is used to retrieve the tokens added during the tokenizer's training. Here's an example of how to use it:



You

LA You

And how would i now use `prepare_for_tokenization`?



ChatGPT

The `prepare_for_tokenization` method in Huggingface's Transformers library allows you to preprocess text before tokenization. This method is useful when you want to tokenize text using the



You

I don't understand the third line



ChatGPT

My apologies for the confusion. Let me clarify the third line of code.

Idea

Chatbots

- Helpful when working and learning new tools, but often cumbersome to use
 - Might need long prompts
 - First results may be unhelpful
 - Multiple follow up prompts can be annoying
- Different LLMs give different results (of different quality)

[f convert_ids_to_tokens](#) <source>

```
( ids: typing.Union[int, typing.List[int]], skip_special_tokens: bool = False ) → str or List[str]
```

Parameters

- **ids** (int or List[int]) — The token id (or token ids) to convert to tokens.
- **skip_special_tokens** (bool, *optional*, defaults to False) — Whether or not to remove special tokens in the decoding.

Returns str or List[str]

The decoded token(s).

Converts a single index or a sequence of indices in a token or a sequence of tokens, using the vocabulary and added tokens.

APIs

- Some API references are very detailed and helpful, others are not

xLLaMa lets you generate code examples for APIs, without having to use any prompts, making it much faster and convenient to get examples.



Who would use xLLaMa?

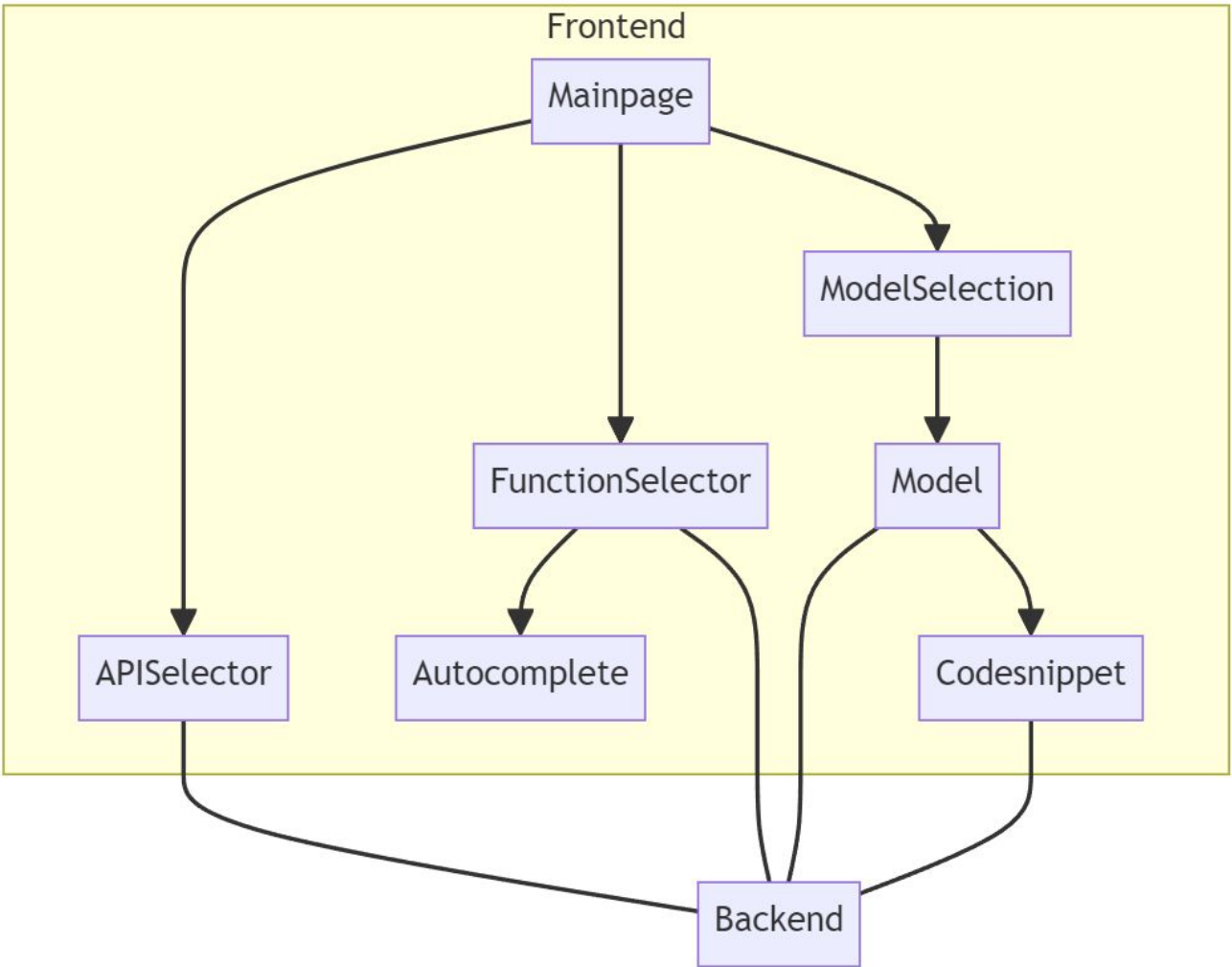
Students and Software Developers who...

- ... are new to one of the frameworks / APIs
- ... find themselves often wondering “wait, how does this work again...?”
- ... want quick results, without having to type a lot to an LLM
- ... want multiple different results to choose from
- ... use multiple different IDEs and source-code editors, or editors without integrated AI assistants


Architecture - Frontend

Vue.js

- Builds on top of standard HTML, CSS, and JavaScript
- Component-based programming model



Architecture - Frontend



API / Framework: ☐ pandas ☒ Hugging Face ☐ numpy

FUNCTION:

tokenizer.batch_decode
☒ codellama ☐ wizardcoder ☐ GPT-3.5 ☐ GPT-4

```

from transformers import AutoTokenizer, AutoModelForSequenceClassification

# Load pre-trained model and tokenizer
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=8)

# Encode a list of sentences as tokenized input ids
sentences = ["This is a sentence.", "This is another sentence."]
encoded_inputs = tokenizer(sentences, return_tensors="pt")

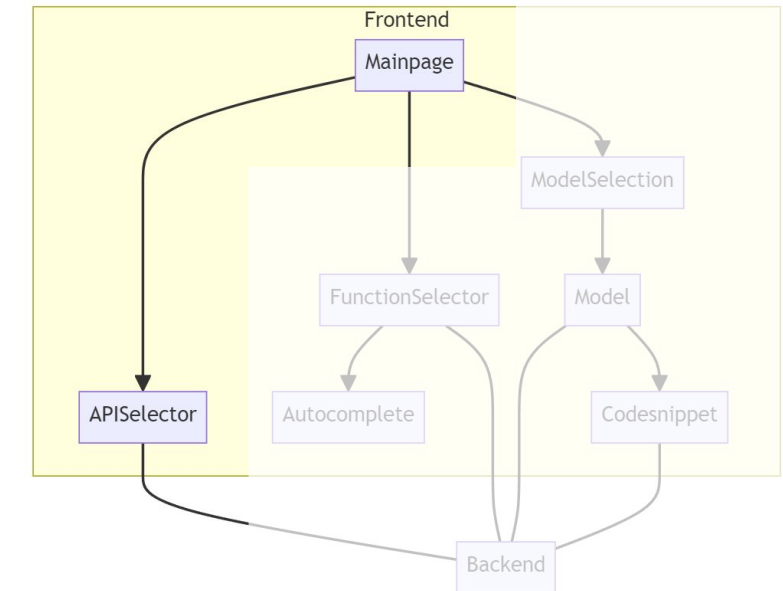
# Decode the encoded input ids into strings
decoded_sentences = model.decode(encoded_inputs["input_ids"], skip_special_tokens=True)
print(decoded_sentences)

```


The APISelector sends a request to the database, asking for the available APIs.

This code is using the Hugging Face Transformers library. Here's a step-by-step explanation:

1. from transformers import AutoTokenizer, AutoModelForSequenceClassification from the Hugging Face Transformers library
2. tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased") assigns it to the variable tokenizer



Architecture - Frontend



API / Framework: 🐼 pandas 🤖 Hugging Face 🔢 numpy

Function: tokenizer.batch_decode tokenizer.convert_ids_to_tokens tokenizer.add_tokens tokenizer.add_special_tokens tokenizer.apply_chat_template

Generate Examples!

tokenizer.batch_decode
🔵 codellama 🔵 wizardcoder 🔵 GPT-3.5 🔵 GPT-4

```

from transformers import AutoTokenizer, AutoModelForSequenceClassification

# Load pre-trained model and tokenizer
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=8)

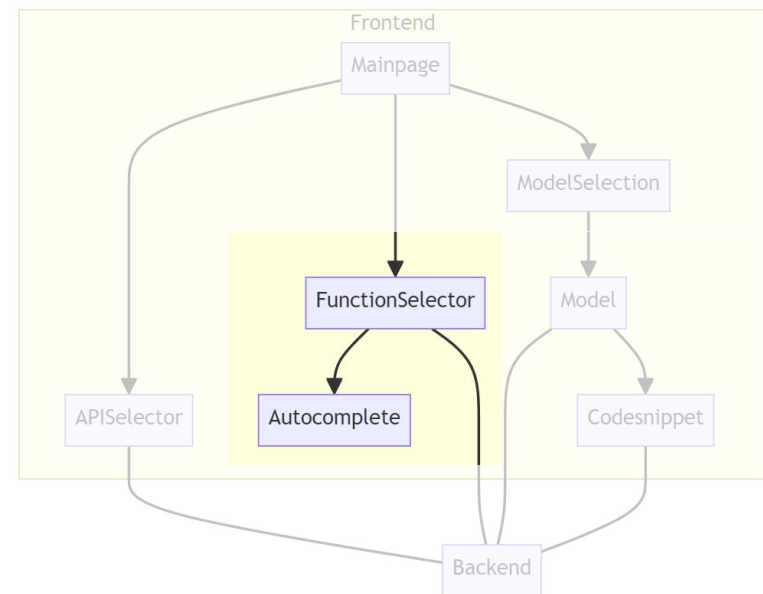
# Encode a list of sentences as tokenized input ids
sentences = ["This is a sentence.", "This is another sentence."]
encoded_inputs = tokenizer(sentences, return_tensors="pt")

# Decode the encoded input ids into strings
decoded_sentences = model.decode(encoded_inputs["input_ids"], skip_special_tokens=True)
print(decoded_sentences)


```

Too long Too short

The FunctionSelector receives the API selected in the previous step and sends a request to the database, asking for the methods available for this API.



Architecture - Frontend



API / Framework: ☐ pandas ☒ Hugging Face ☐ numpy

Function:

tokenizer.batch_decode
☒ codellama ☐ wizardcoder ☐ GPT-3.5 ☐ GPT-4

```

from transformers import AutoTokenizer, AutoModelForSequenceClassification

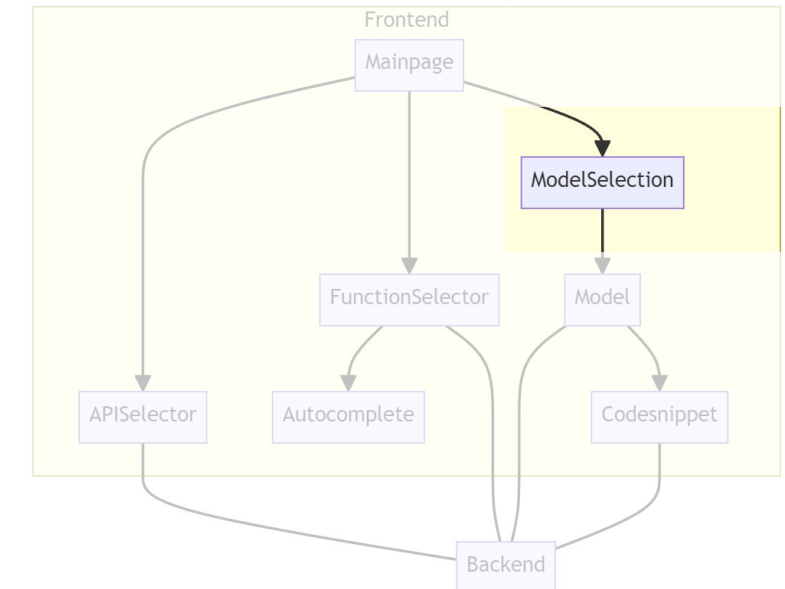
# Load pre-trained model and tokenizer
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=8)

# Encode a list of sentences as tokenized input ids
sentences = ["This is a sentence.", "This is another sentence."]
encoded_inputs = tokenizer(sentences, return_tensors="pt")


# Decode the encoded input ids into strings
decoded_sentences = model.decode(encoded_inputs["input_ids"], skip_special_tokens=True)
print(decoded_sentences)

```

The ModelSelection passes the request to generate an example to the different Model instances.



Architecture - Frontend



API / Framework: ☐ pandas ☒ Hugging Face ☐ numpy

Function:

tokenizer.batch_decode

☒ codellama ☐ wizardcoder ☐ GPT-3.5 ☐ GPT-4

```
from transformers import AutoTokenizer, AutoModelForSequenceClassification

# Load pre-trained model and tokenizer
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=8)

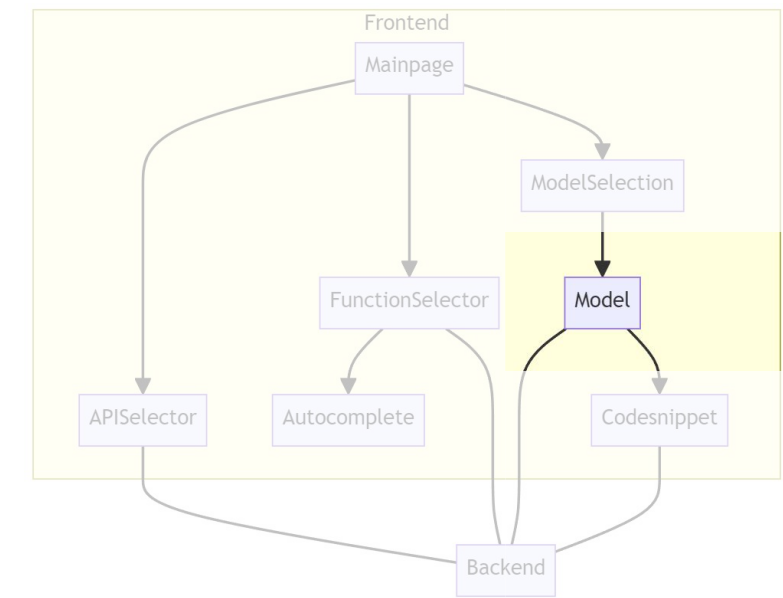
# Encode a list of sentences as tokenized input ids
sentences = ["This is a sentence.", "This is another sentence."]
encoded_inputs = tokenizer(sentences, return_tensors="pt")

# Decode the encoded input ids into strings
decoded_sentences = model.decode(encoded_inputs["input_ids"], skip_special_tokens=True)
print(decoded_sentences)
```


A Model instance sends the request to generate an example to the backend and receives the datastream.

This code is using the Hugging F
Here's a step-by-step explanatio

1. from transformers import Au
from the Hugging Face Tra
2. tokenizer = AutoTokenizer.f
assigns it to the variable to



Architecture - Frontend



API / Framework: ☐ pandas ☒ Hugging Face ☐ numpy

Function:

tokenizer.batch_decode
☒ codellama ☐ wizardcoder ☐ GPT-3.5 ☐ GPT-4

```

from transformers import AutoTokenizer, AutoModelForSequenceClassification

# Load pre-trained model and tokenizer
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=8)

# Encode a list of sentences as tokenized input ids
sentences = ["This is a sentence.", "This is another sentence."]
encoded_inputs = tokenizer(sentences, return_tensors="pt")

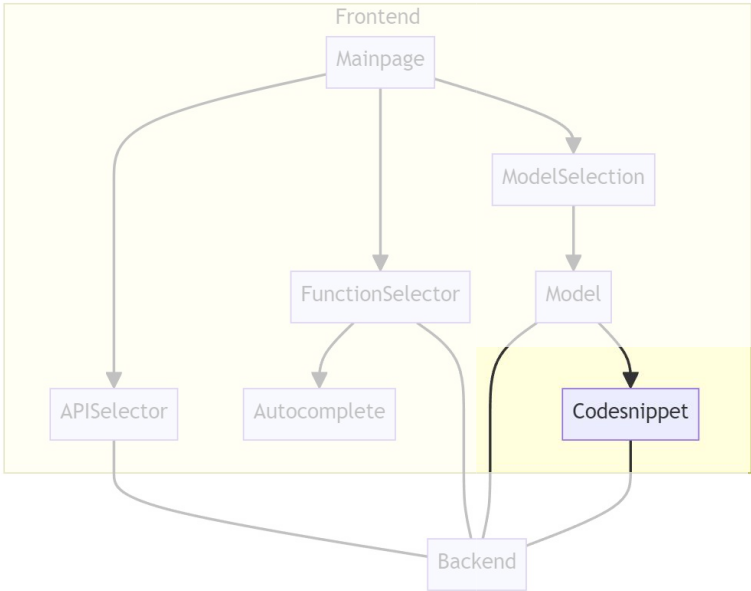
# Decode the encoded input ids into strings
decoded_sentences = model.decode(encoded_inputs["input_ids"], skip_special_tokens=True)
print(decoded_sentences)

```

This code is using the Hugging Face Transformers library. Here's a step-by-step explanation:

1. from transformers import AutoTokenizer, AutoModelForSequenceClassification: This line imports the necessary classes from the Hugging Face Transformers library.
2. tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased"):

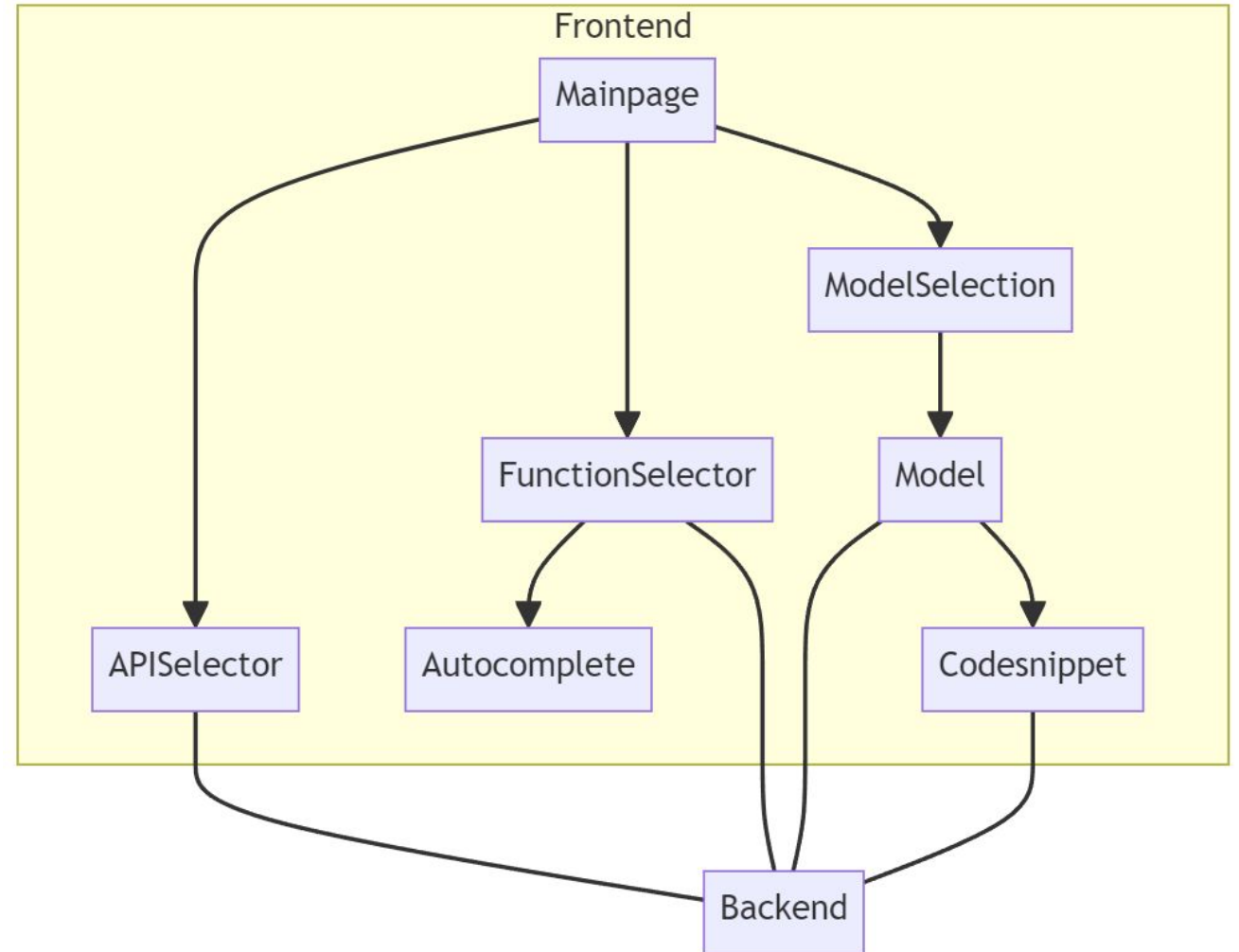
A Codesnippet is a part of a generated example, which only consists of code. Syntaxhighlighting is applied in the backend before it's sent back to the frontend.



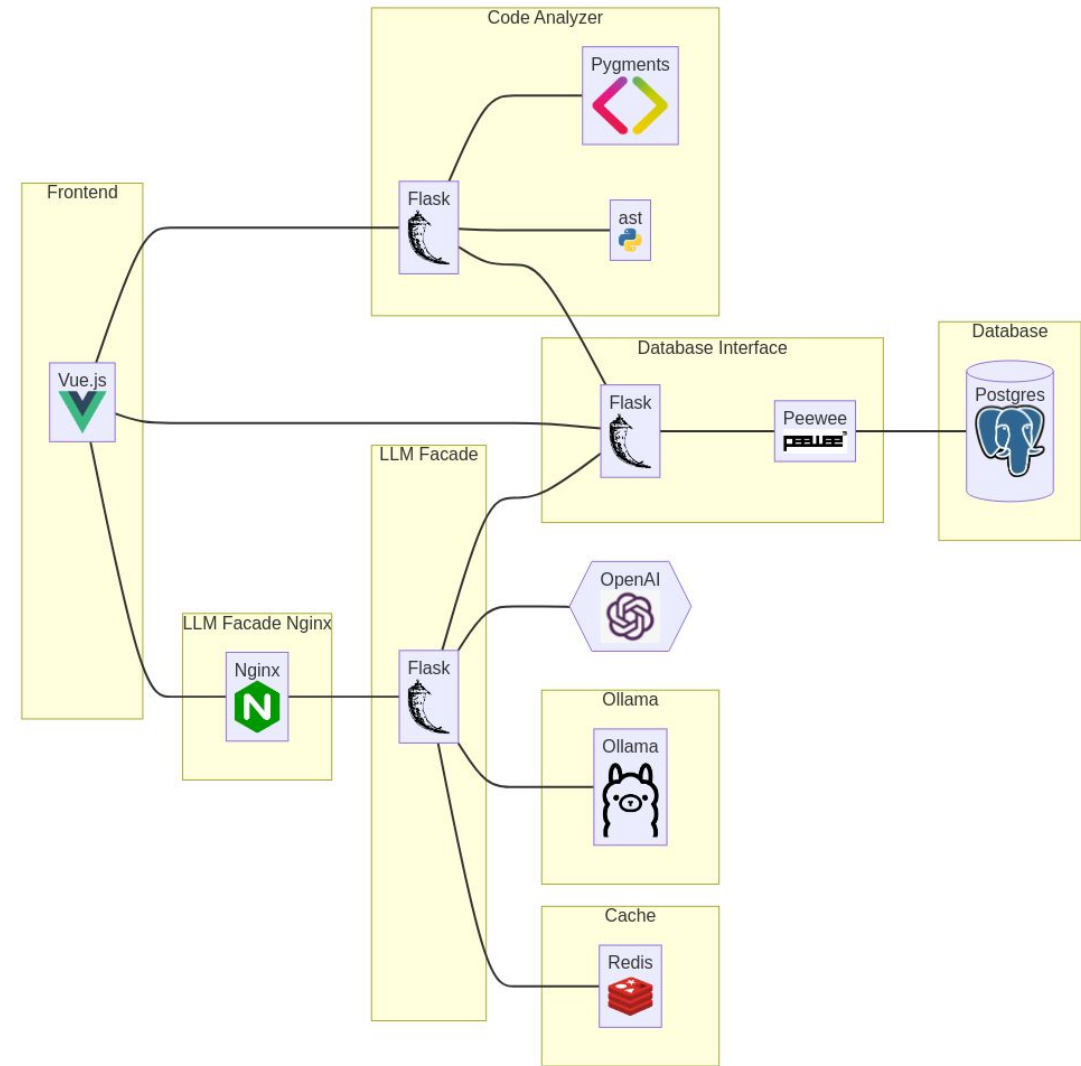
Architecture - Frontend

Vue.js

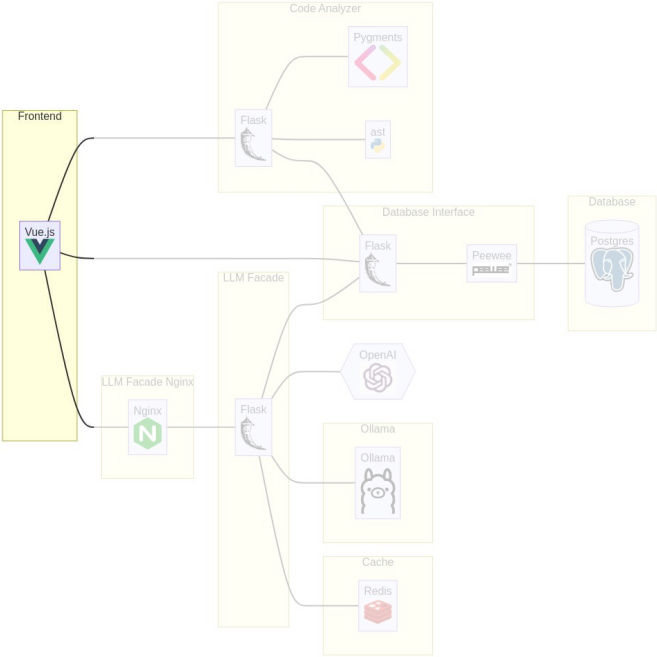
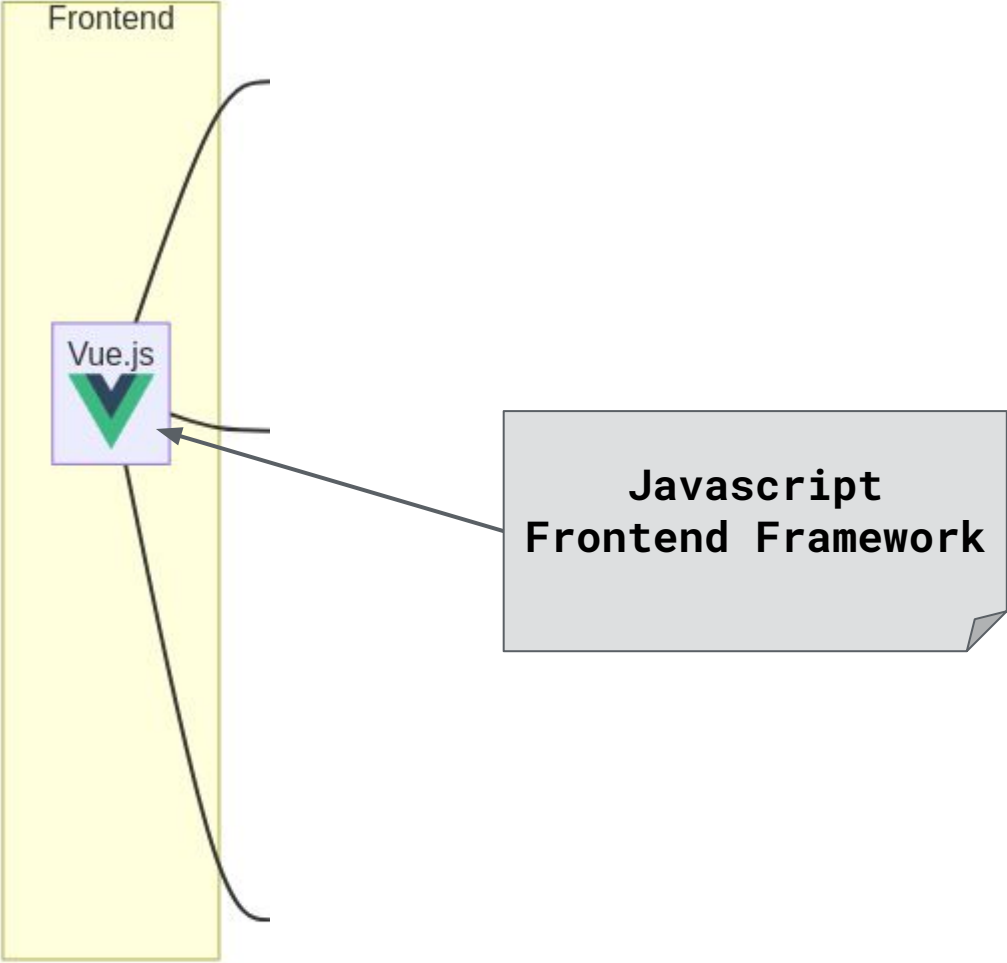
- Builds on top of standard HTML, CSS, and JavaScript
- Component-based programming model



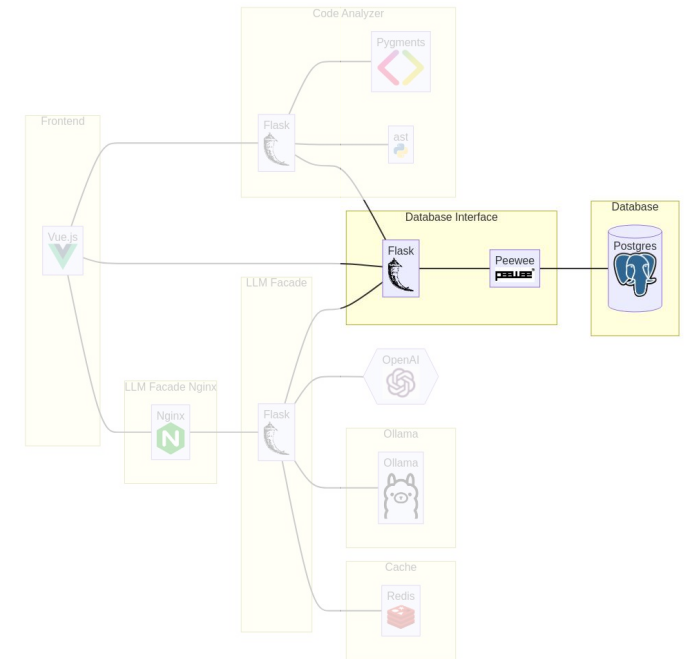
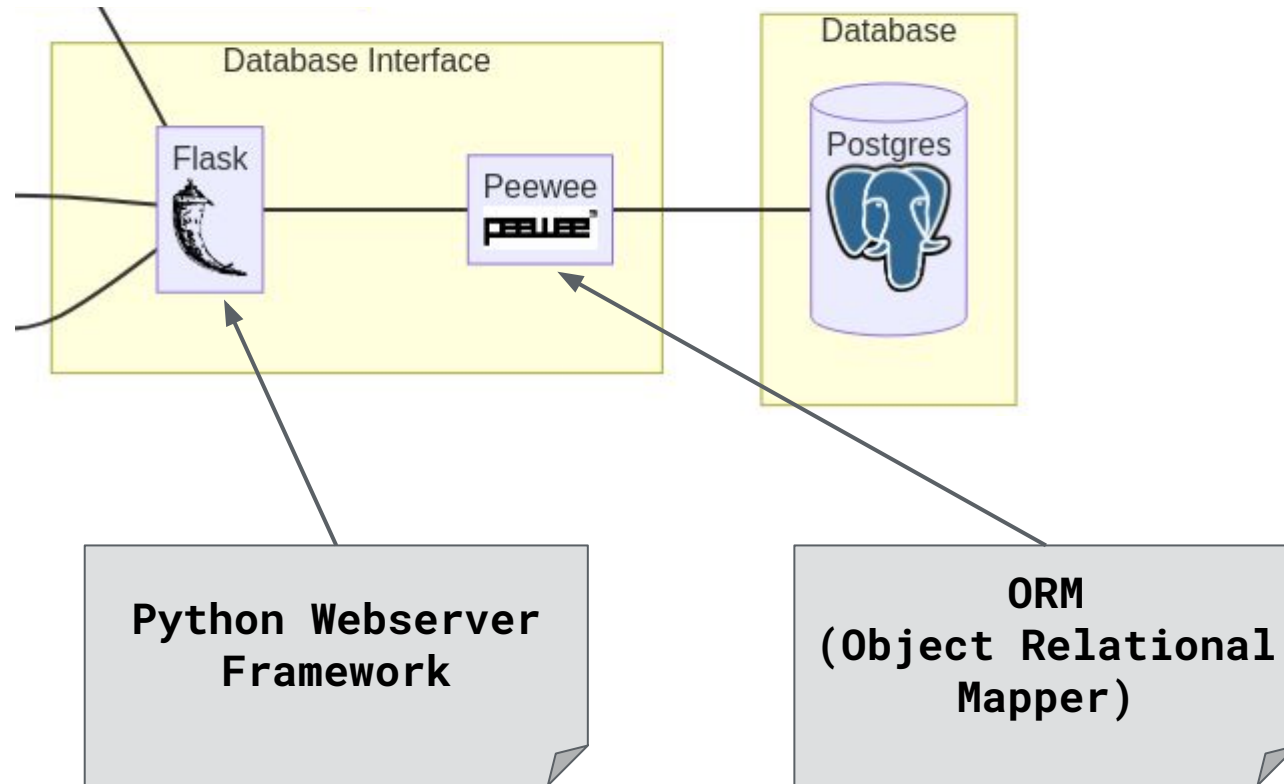
Architecture - Backend



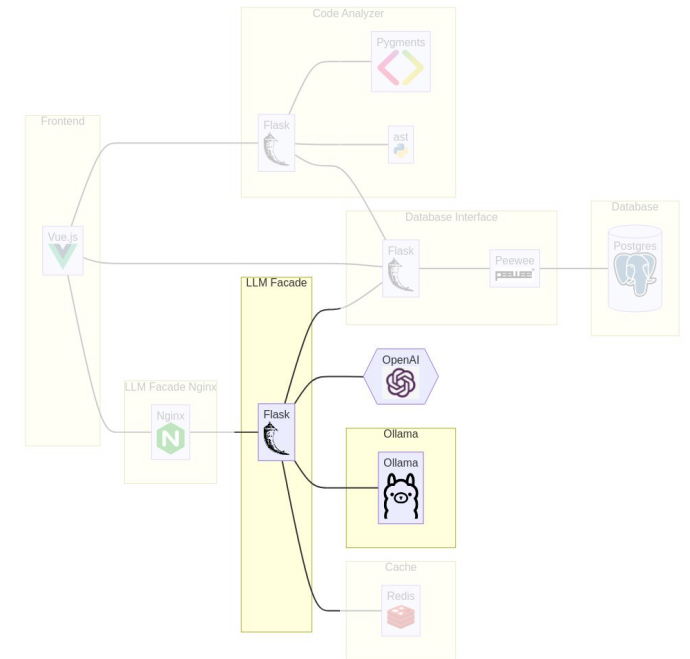
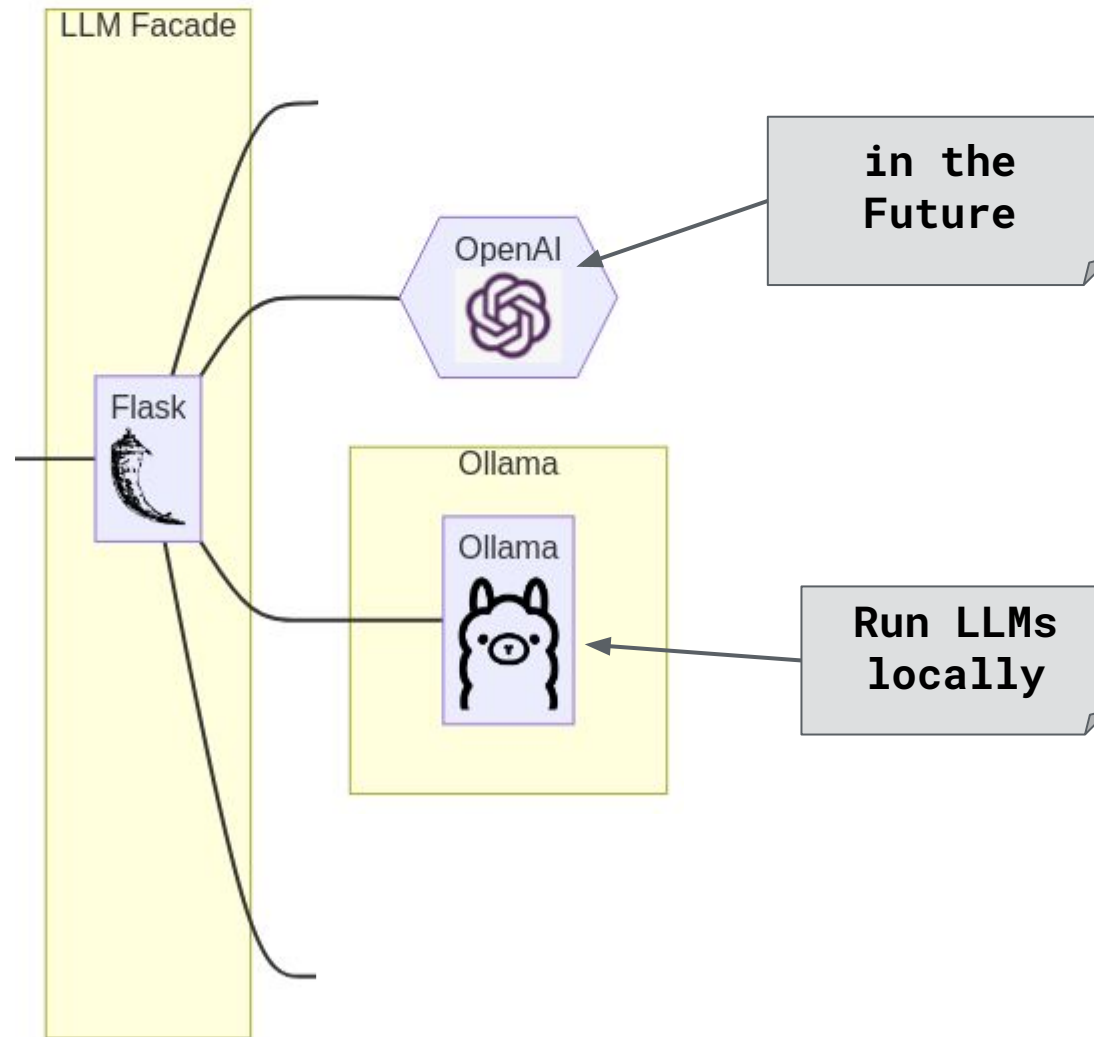
Architecture - Backend



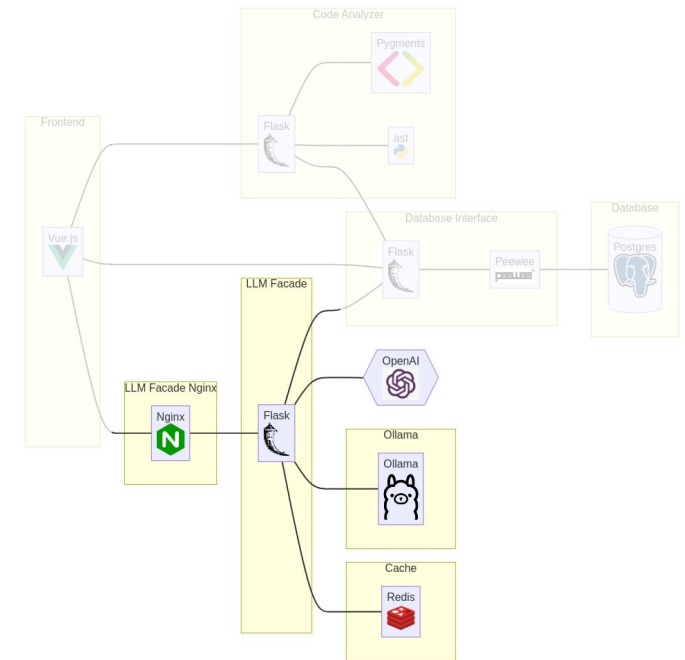
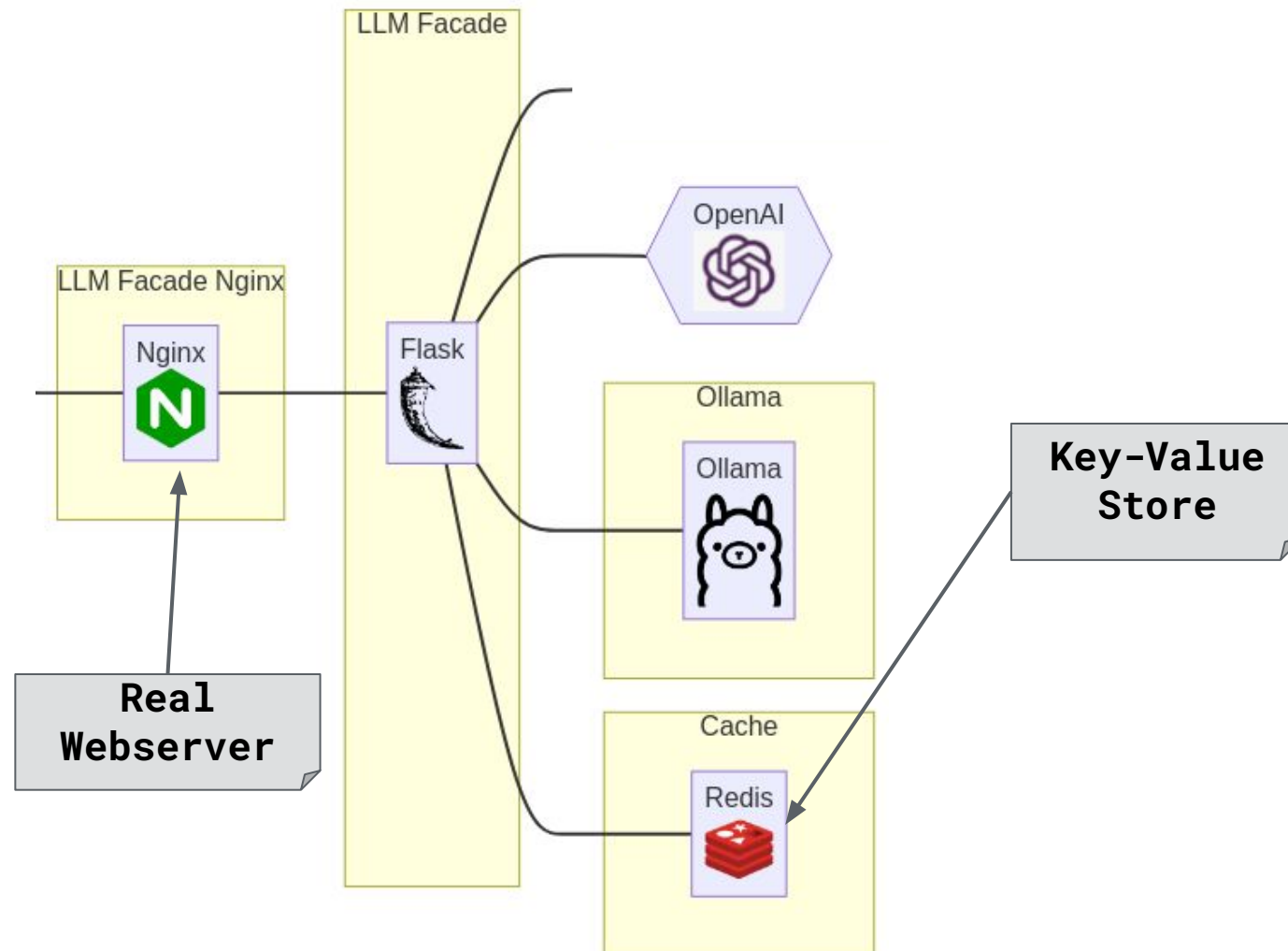
Architecture - Backend



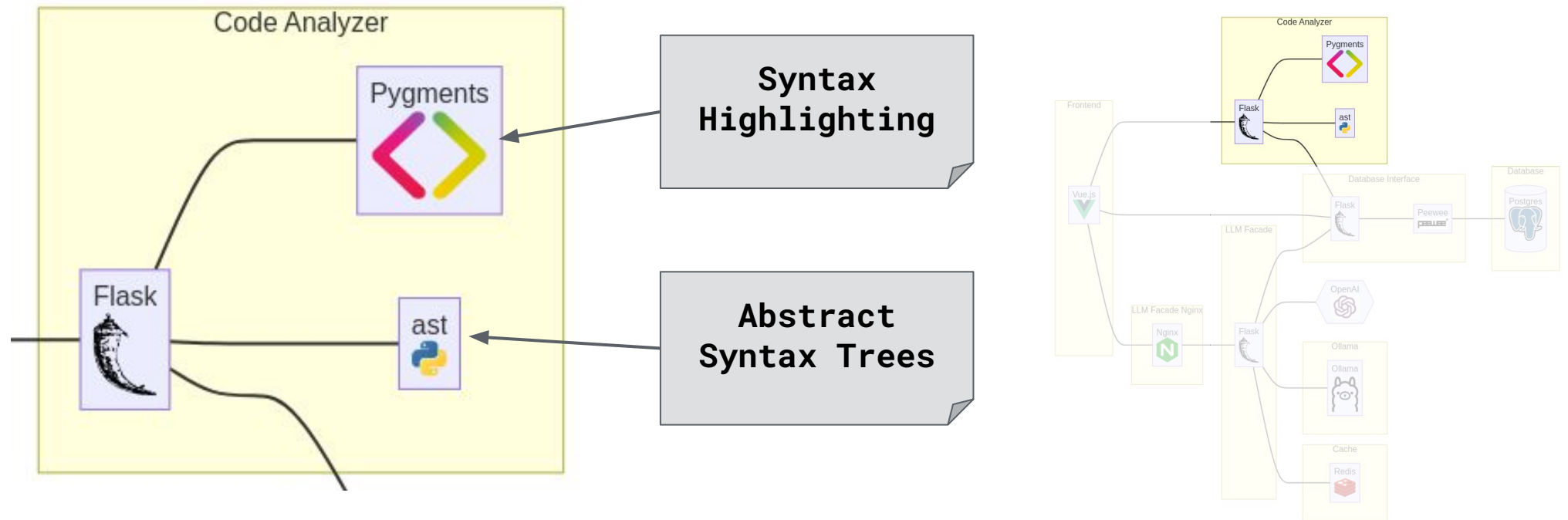
Architecture - Backend



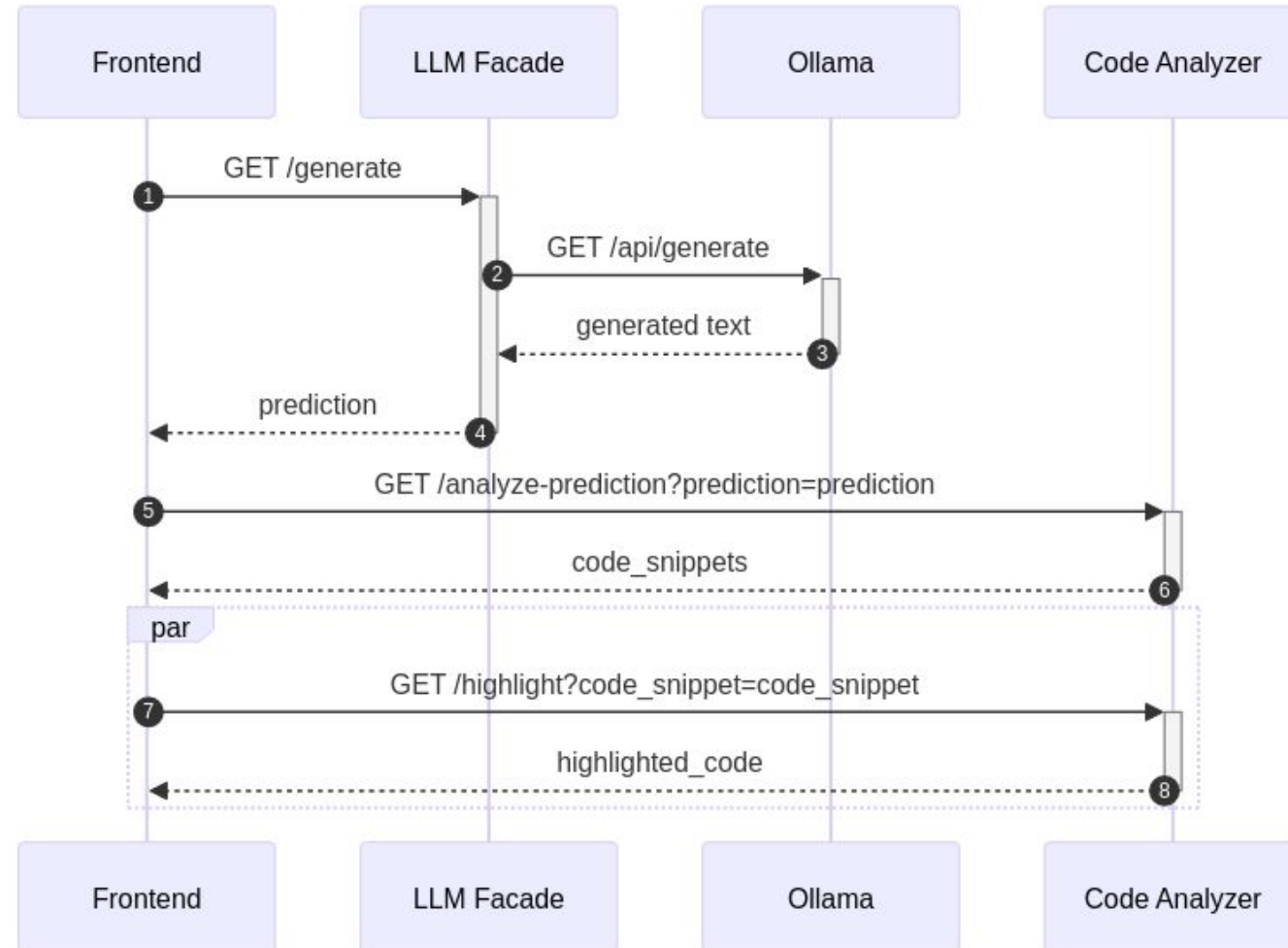
Architecture - Backend



Architecture - Backend




Example Generation




Example Generation

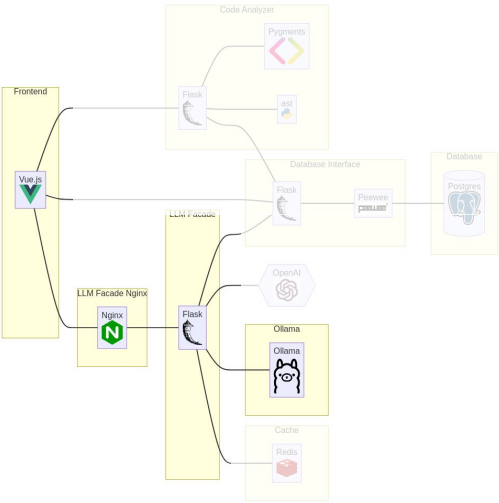
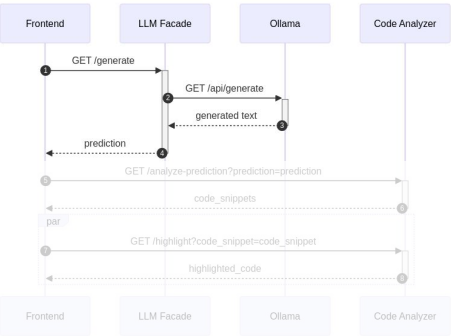
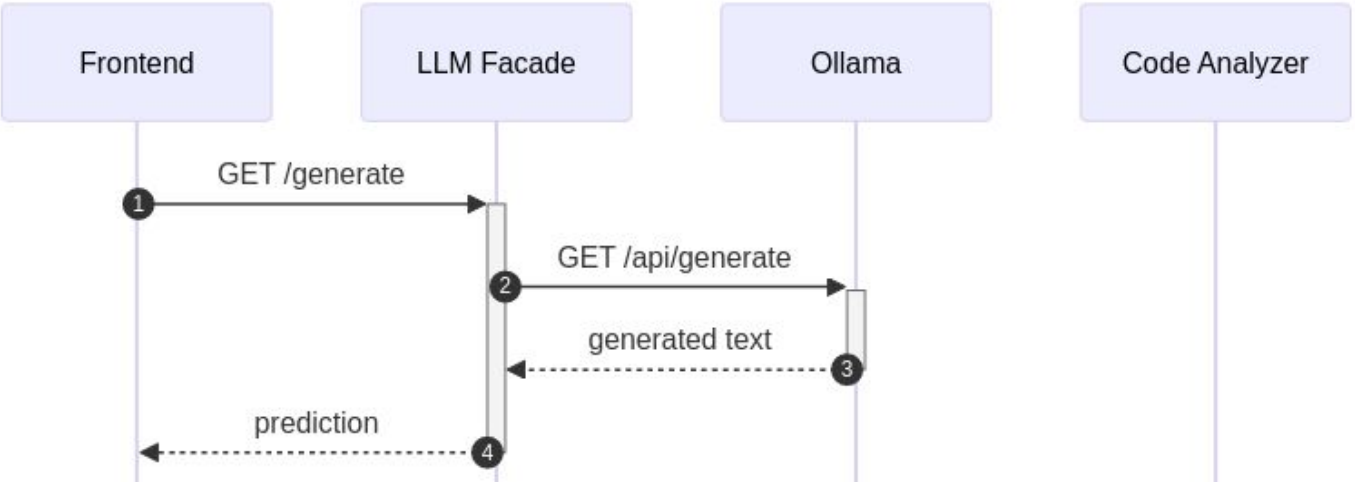
DataFrame.at

☒ codellama ☐ wizardcoder ☐ GPT-3.5 ☐ GPT-4

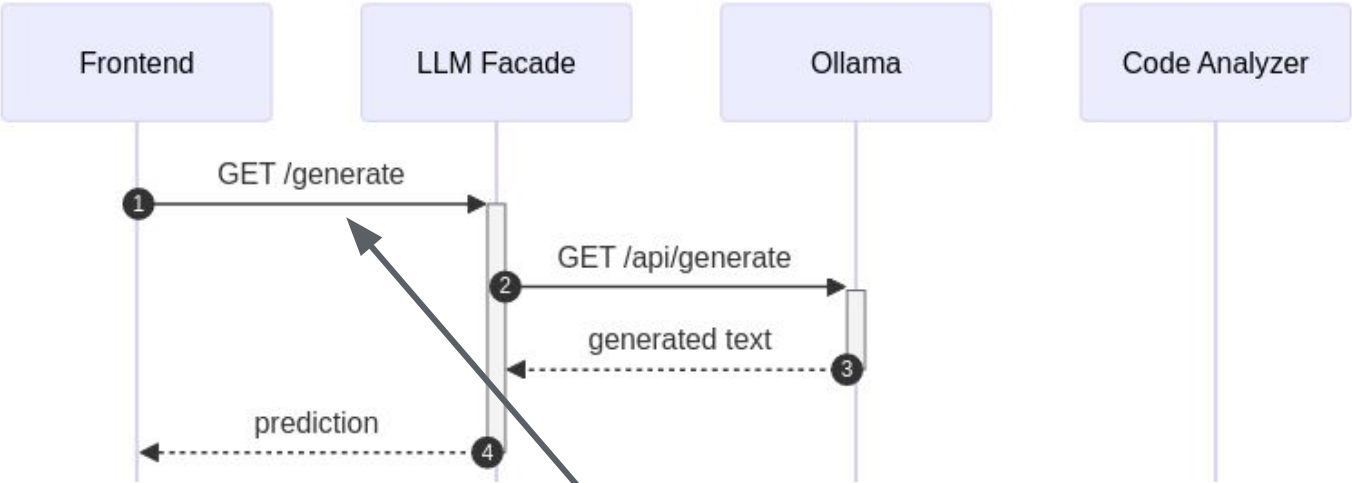




Example Generation

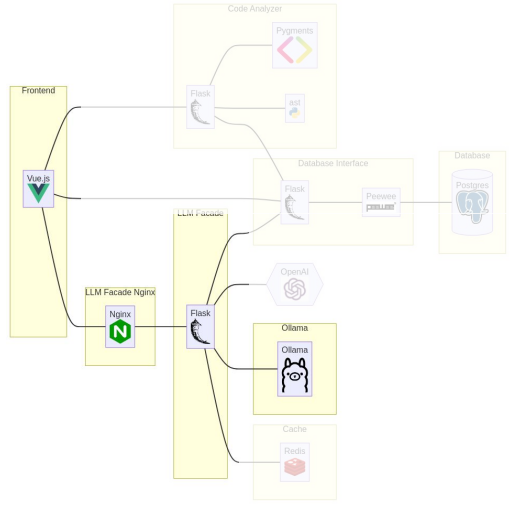
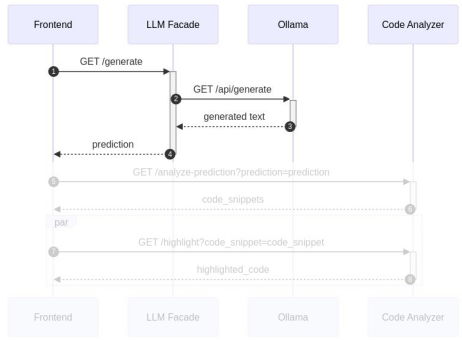


Example Generation

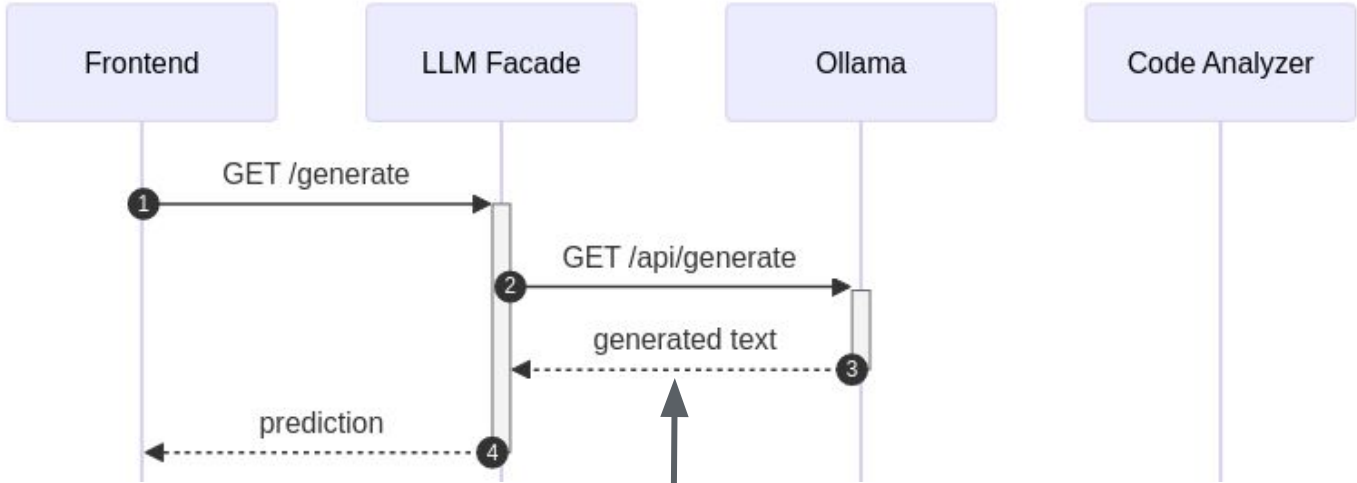


```

{
  model: "codellama:7b-instruct",
  function: "DataFrame.at",
  prompt: "You are a ..."
}
  
```



Example Generation

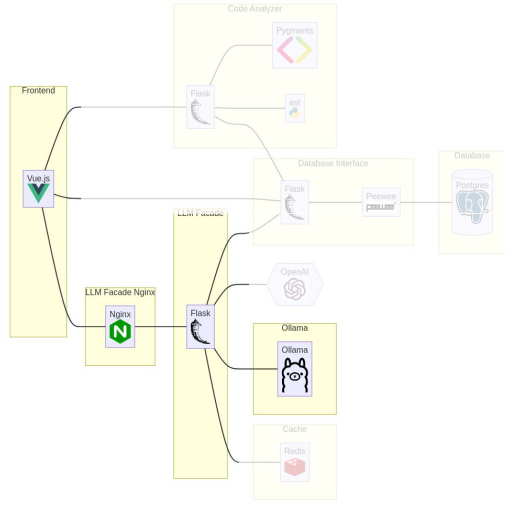
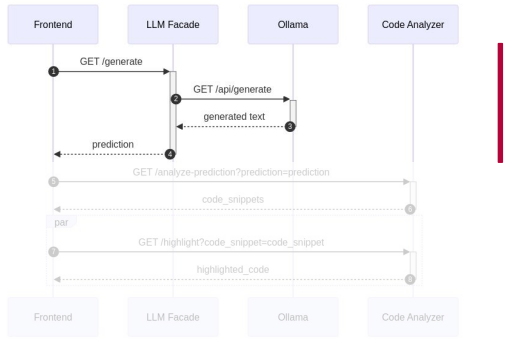


```

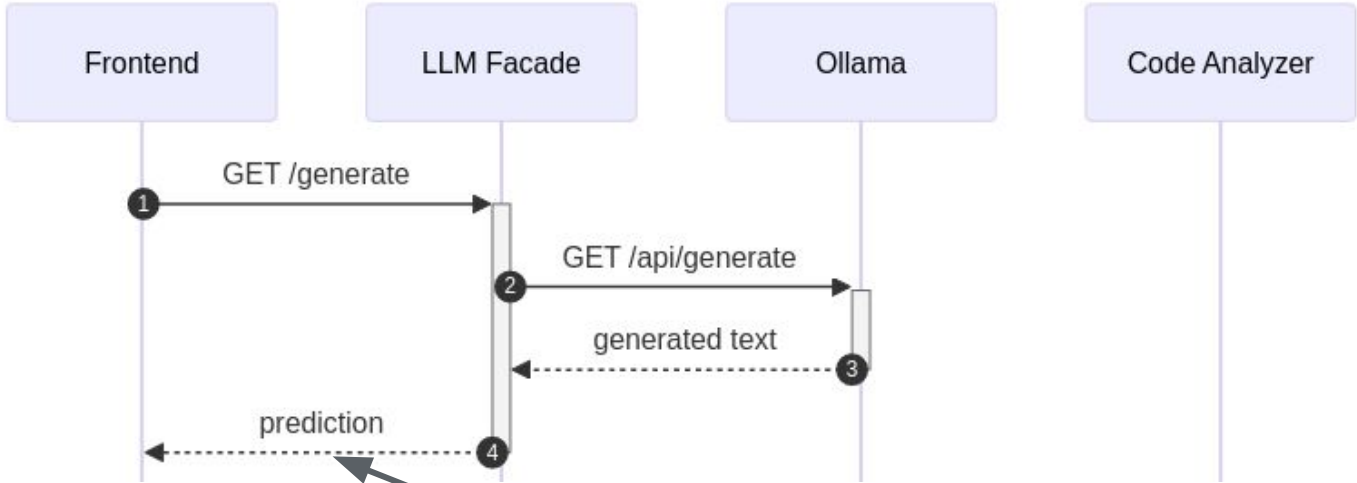
import pandas as pd

# create a sample DataFrame
df = pd.DataFrame({'A': [0, 2, 3],
                   'B': [4, 1, 6], 'C': [3, 2, 1]},
                  index=[4, 5, 6])

...
    
```

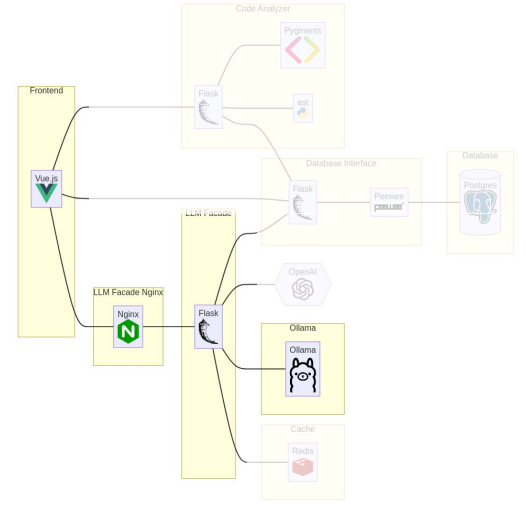
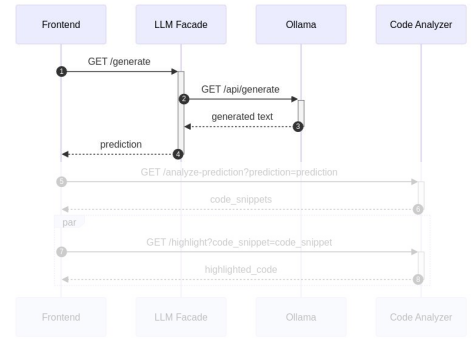


Example Generation



```

{
  text: "import pandas as pd ...",
  model: "codellama:7b-instruct",
  token_amount: 456
}
  
```



Example Generation

DataFrame.at

☒ codellama ☐ wizardcoder ☐ GPT-3.5 ☐ GPT-4

```
```python
import pandas as pd

create a sample DataFrame
df = pd.DataFrame({'A': [0, 2, 3], 'B': [4, 1, 6], 'C': [3, 2, 1]}, index=[4, 5, 6])

use the at property to access a single value in the DataFrame
print(df.at[4, 'B']) # output: 2

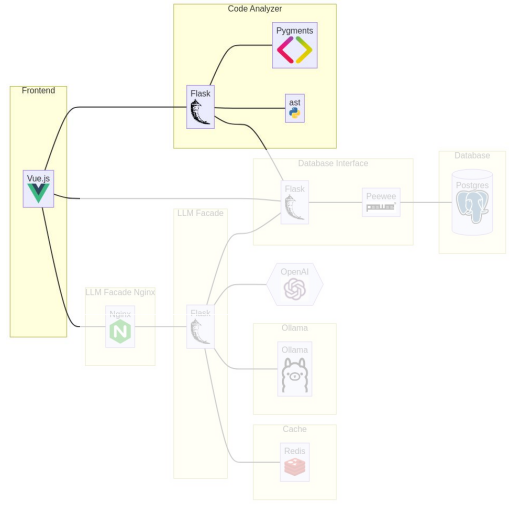
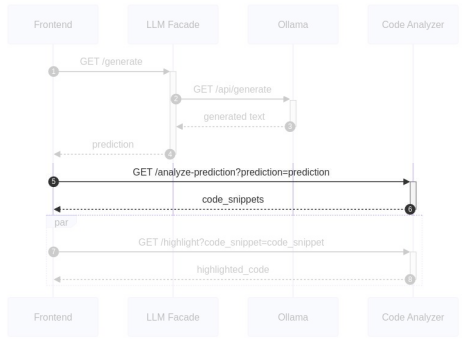
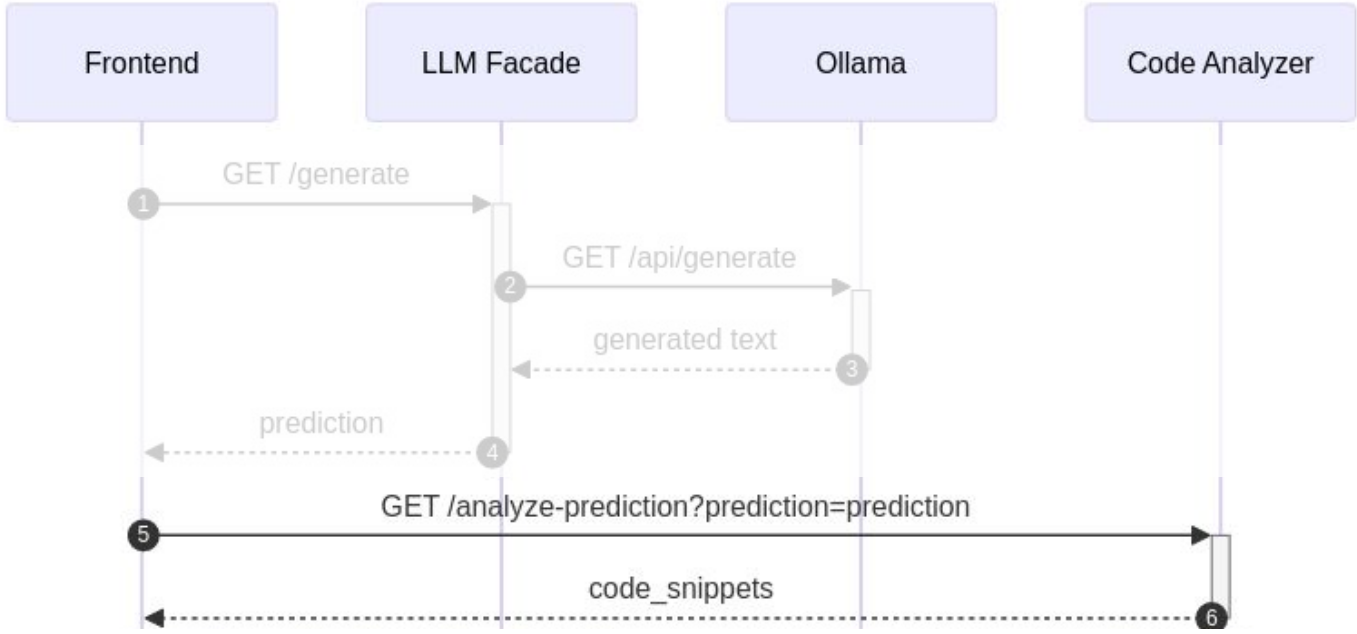
set a new value for a row/column pair using the at property
df.at[4, 'B'] = 10

use the at property to get a single value within a Series
print(df.loc[5].at['B']) # output: 4
```
```

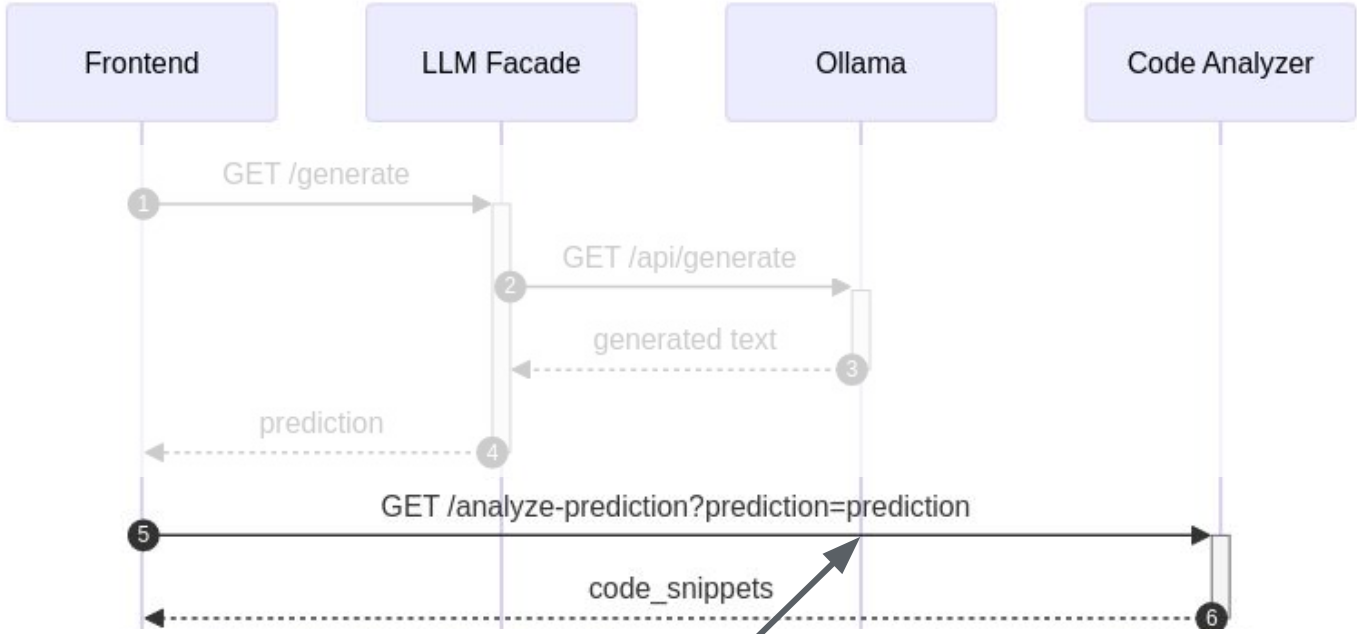
Too long

Too short

Example Generation

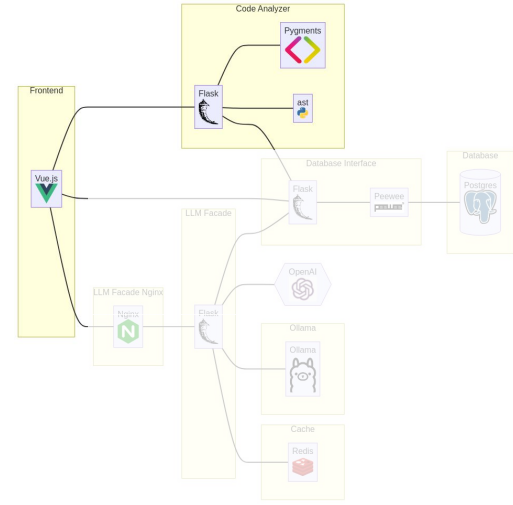
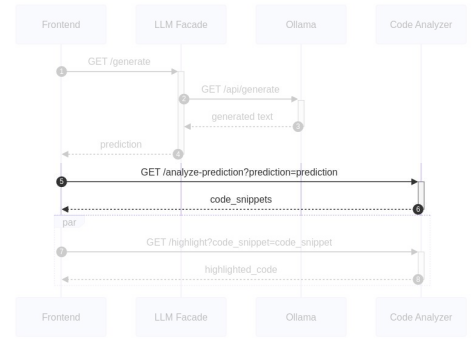


Example Generation

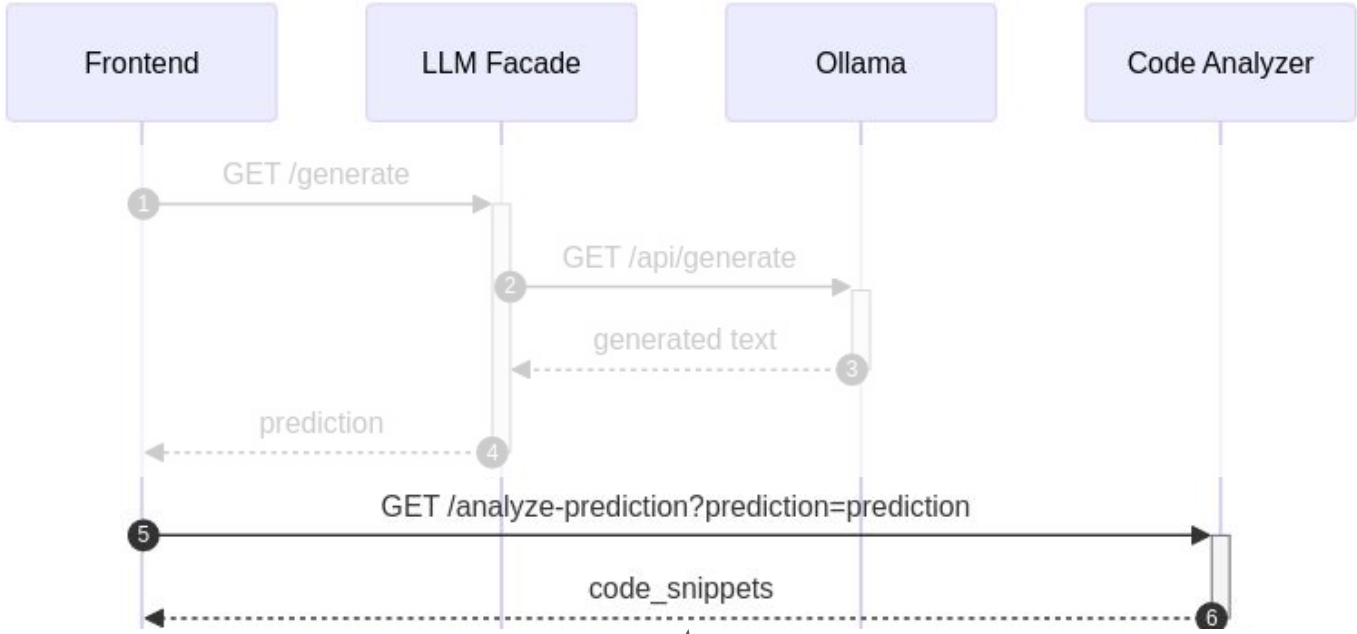


```

{
  text: "import pandas as pd ...",
  model: "codellama:7b-instruct",
  token_amount: 456
}
  
```

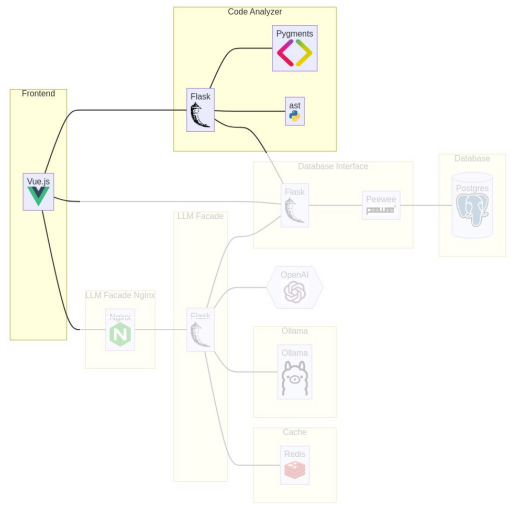
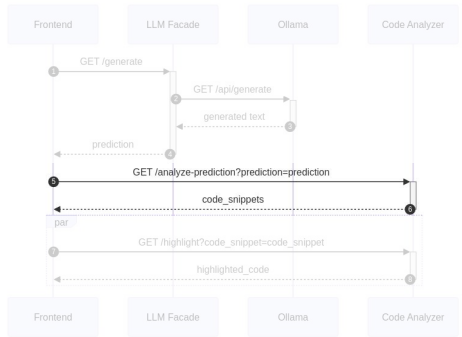


Example Generation

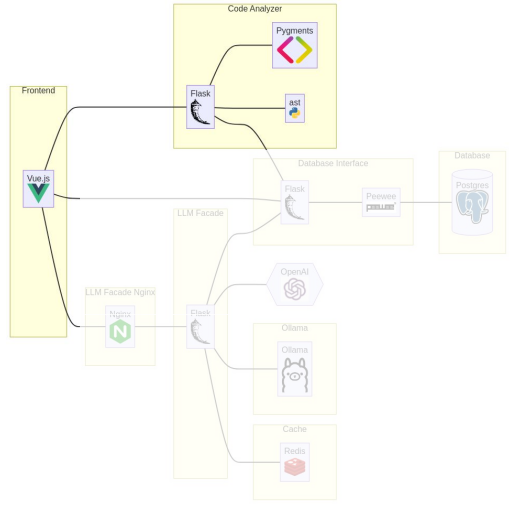
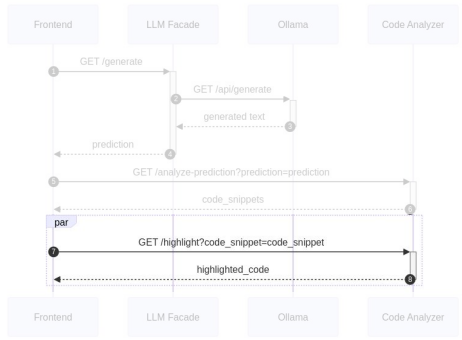
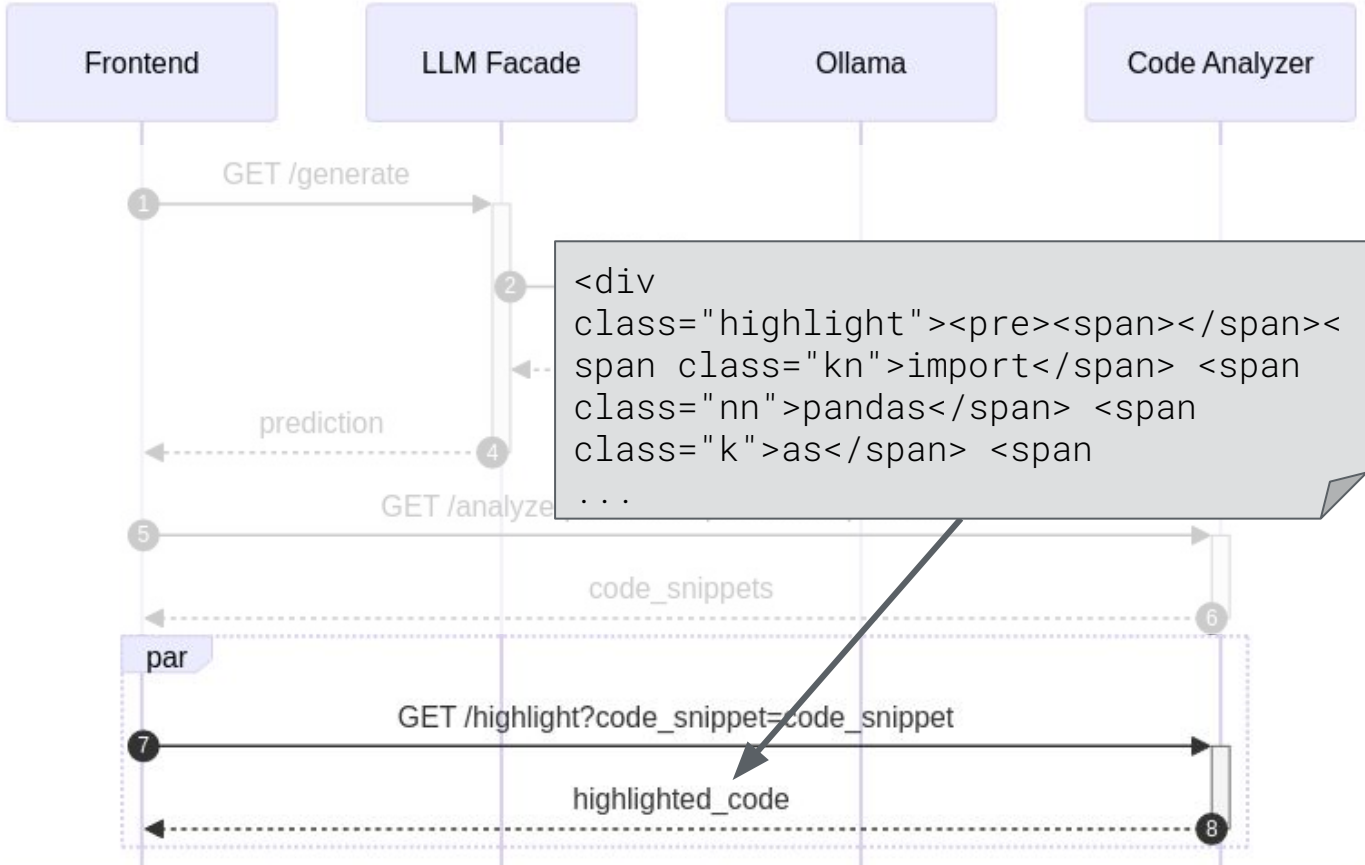


```

[
  {
    text: "import pandas as pd ..."
  }
]
  
```





Example Generation



Example Generation

DataFrame.at

☒ codellama ☐ wizardcoder ☐ GPT-3.5 ☐ GPT-4



```
import pandas as pd

# create a sample DataFrame
df = pd.DataFrame({'A': [0, 2, 3], 'B': [4, 1, 6], 'C': [3, 2, 1]}, index=[4, 5, 6])

# use the at property to access a single value in the DataFrame
print(df.at[4, 'B']) # output: 2

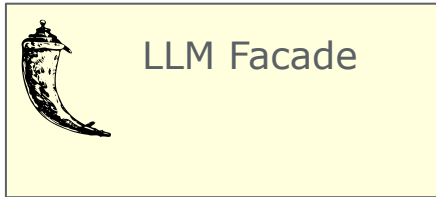
# set a new value for a row/column pair using the at property
df.at[4, 'B'] = 10

# use the at property to get a single value within a Series
print(df.loc[5].at['B']) # output: 4
```

Too long

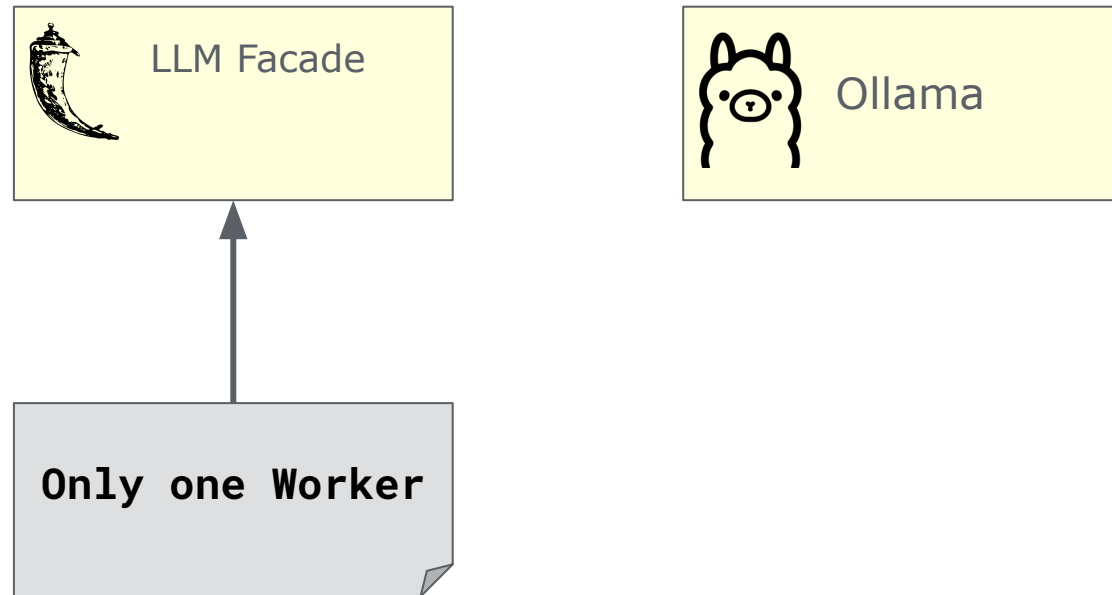
Too short

Simultaneous Example Generation

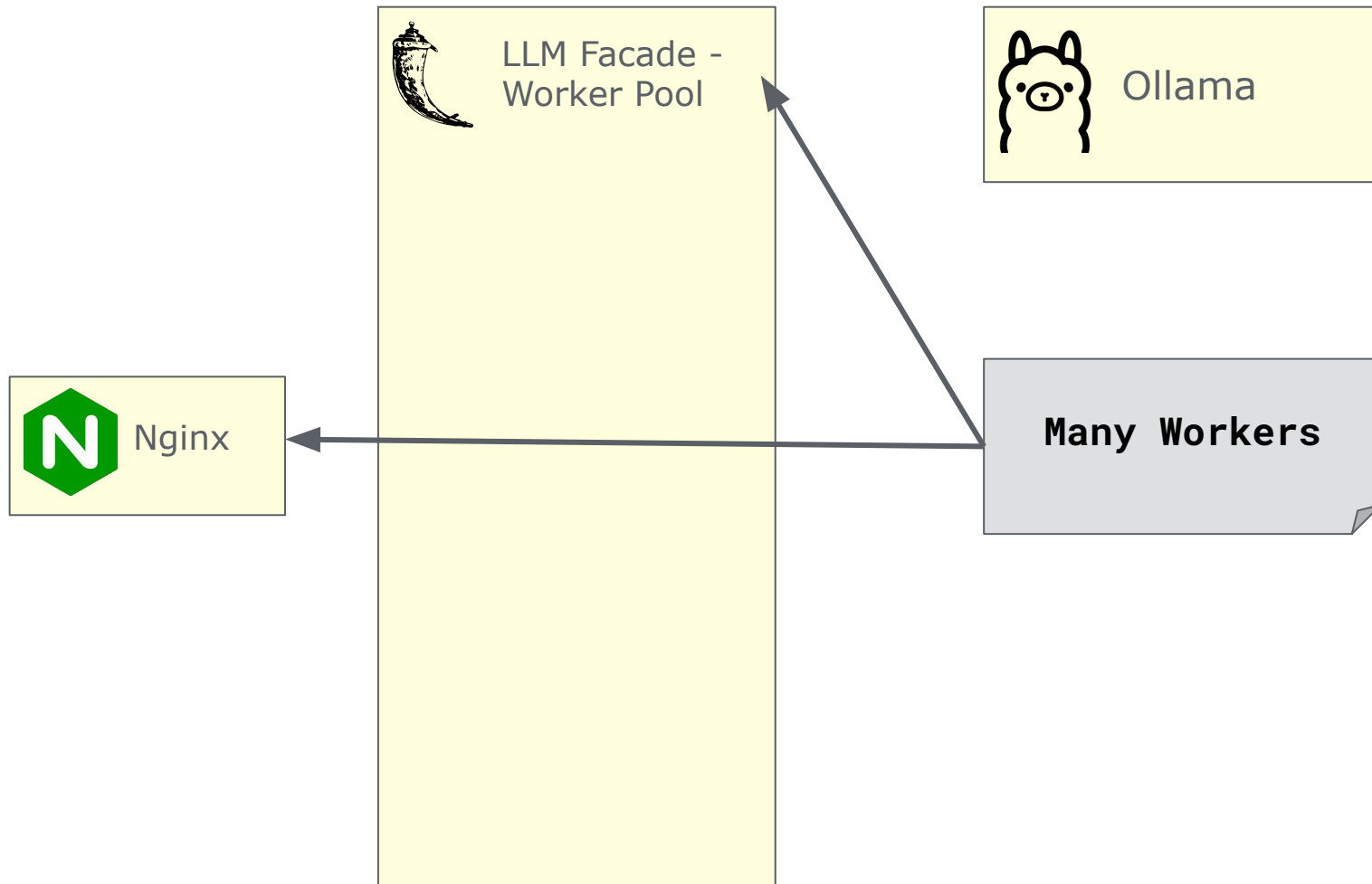


Can it handle multiple requests?

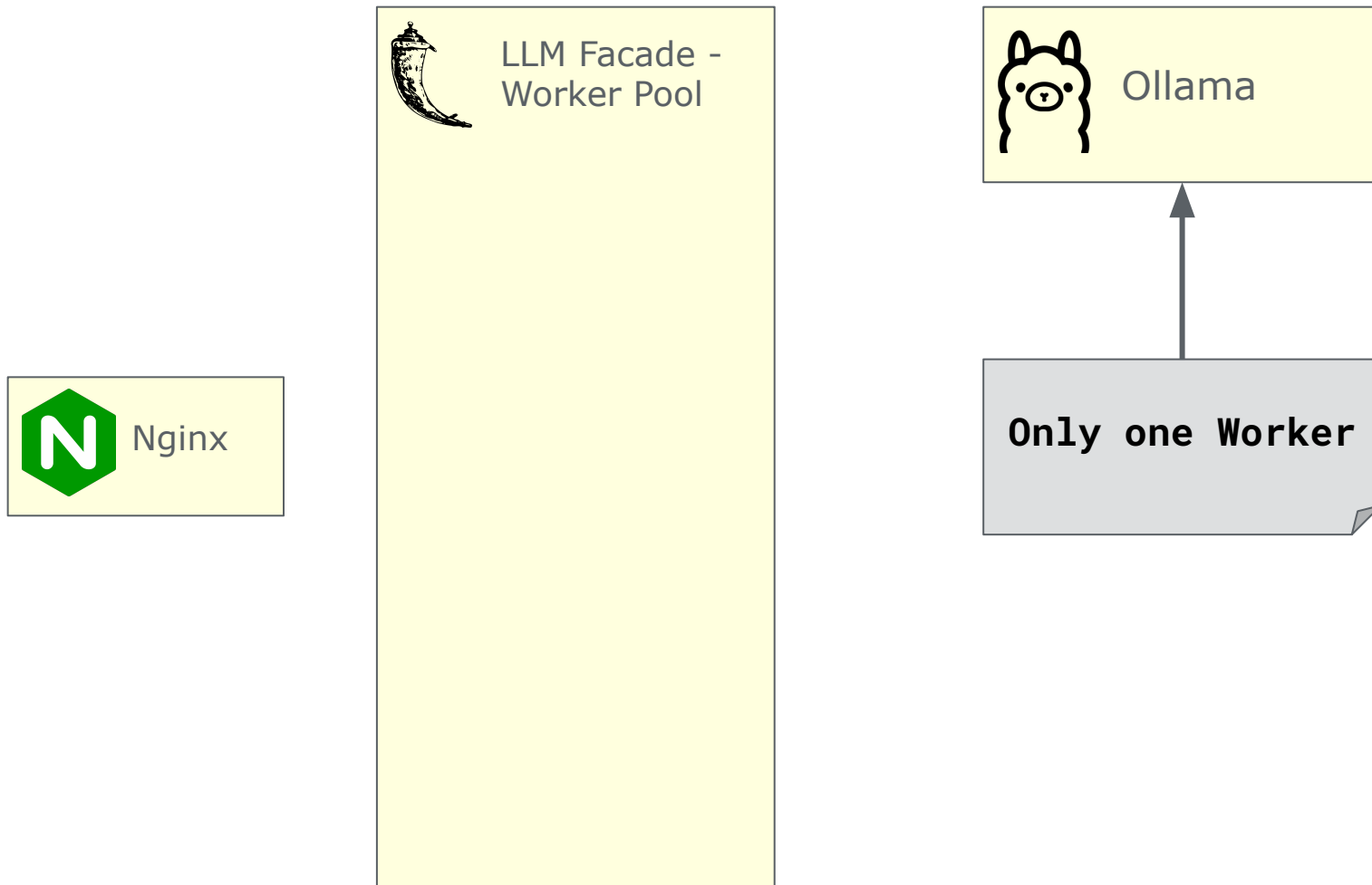
Simultaneous Example Generation



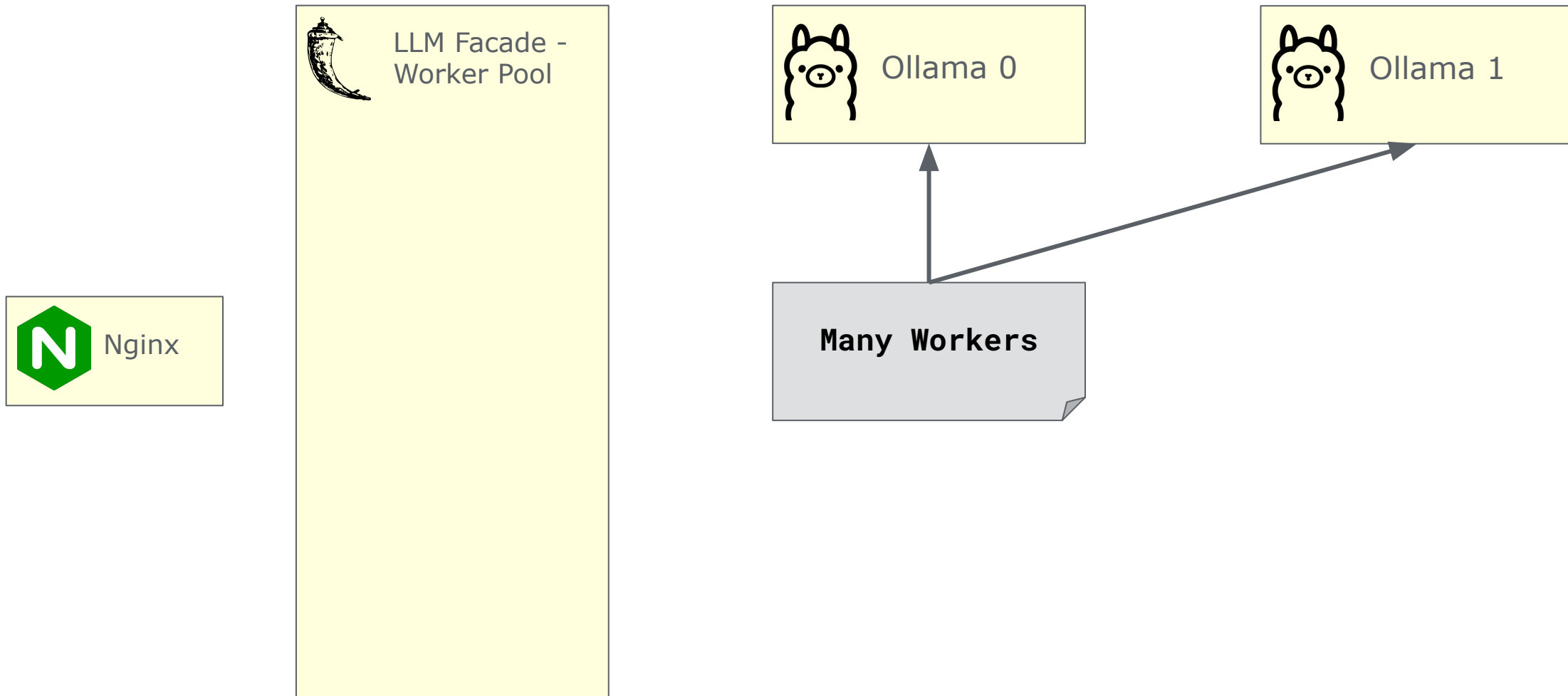
Simultaneous Example Generation



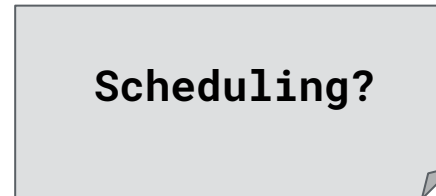
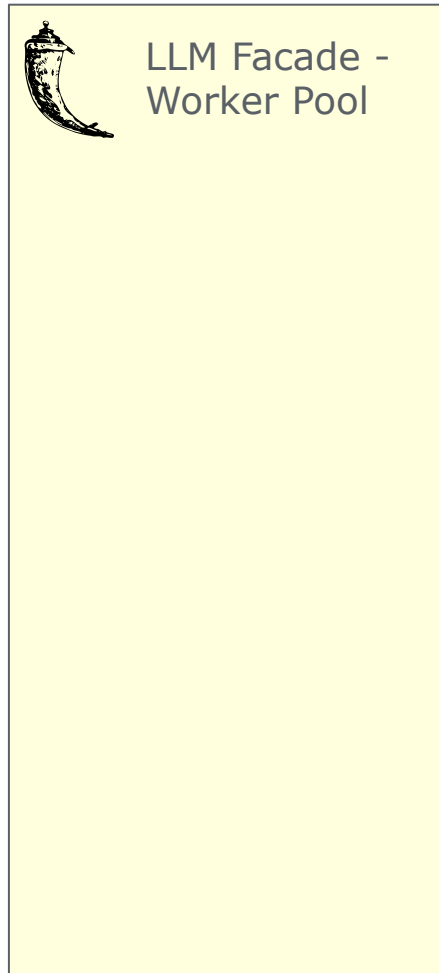
Simultaneous Example Generation



Simultaneous Example Generation



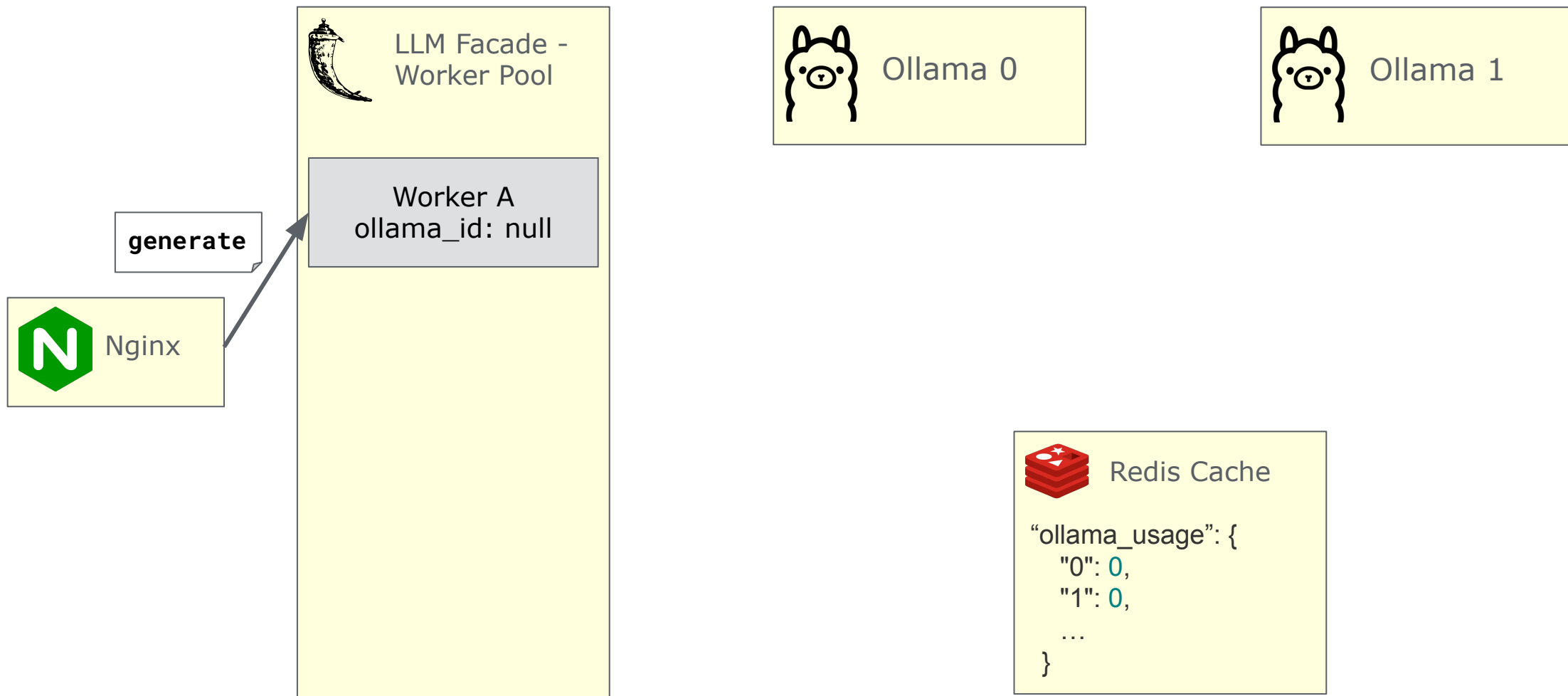
Simultaneous Example Generation



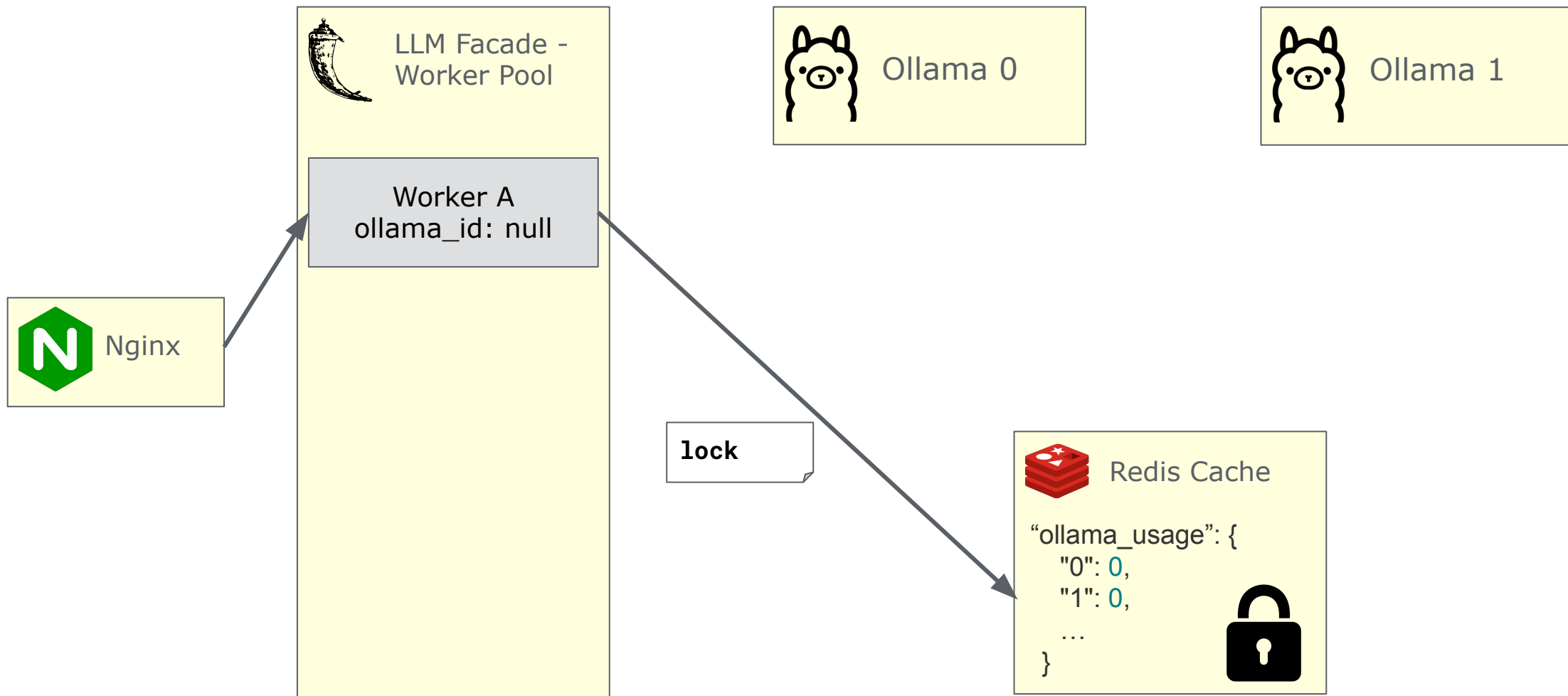
Simultaneous Example Generation



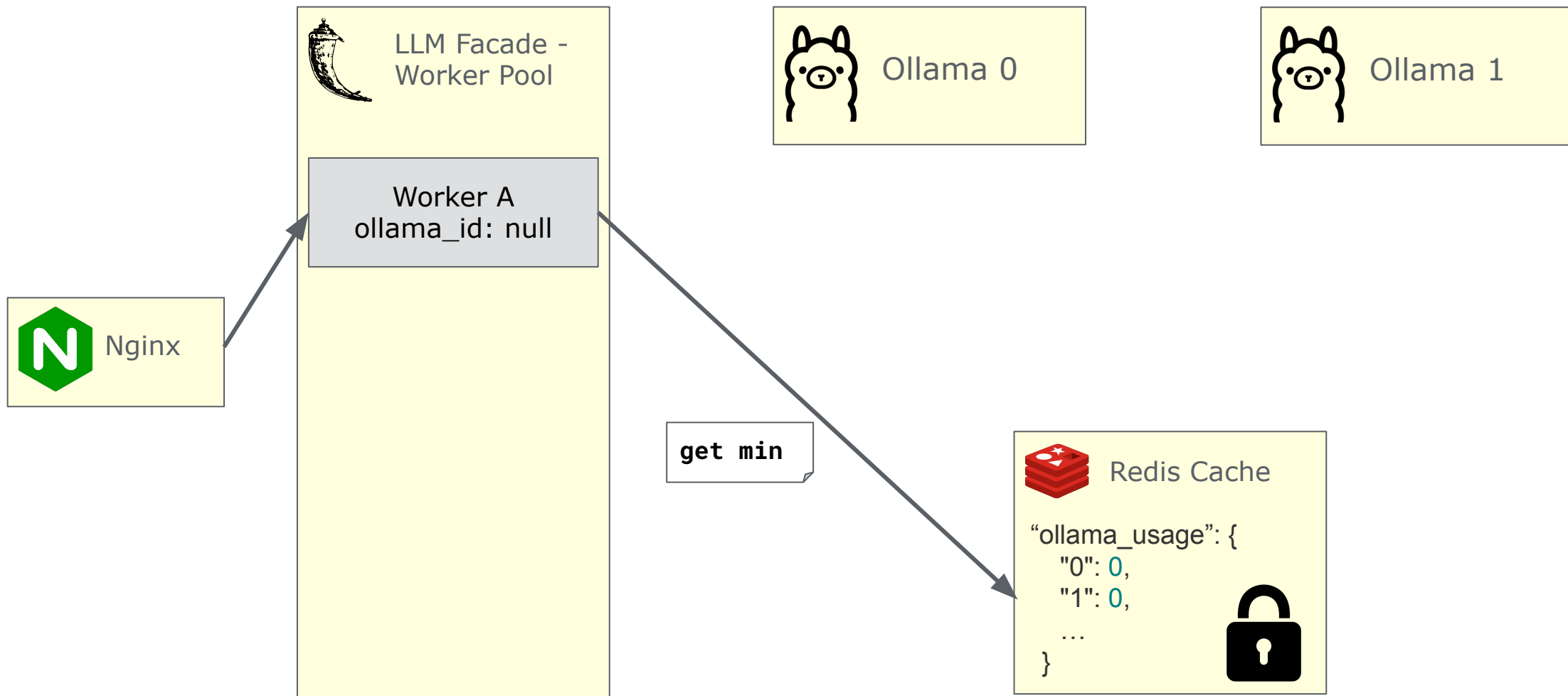
Simultaneous Example Generation



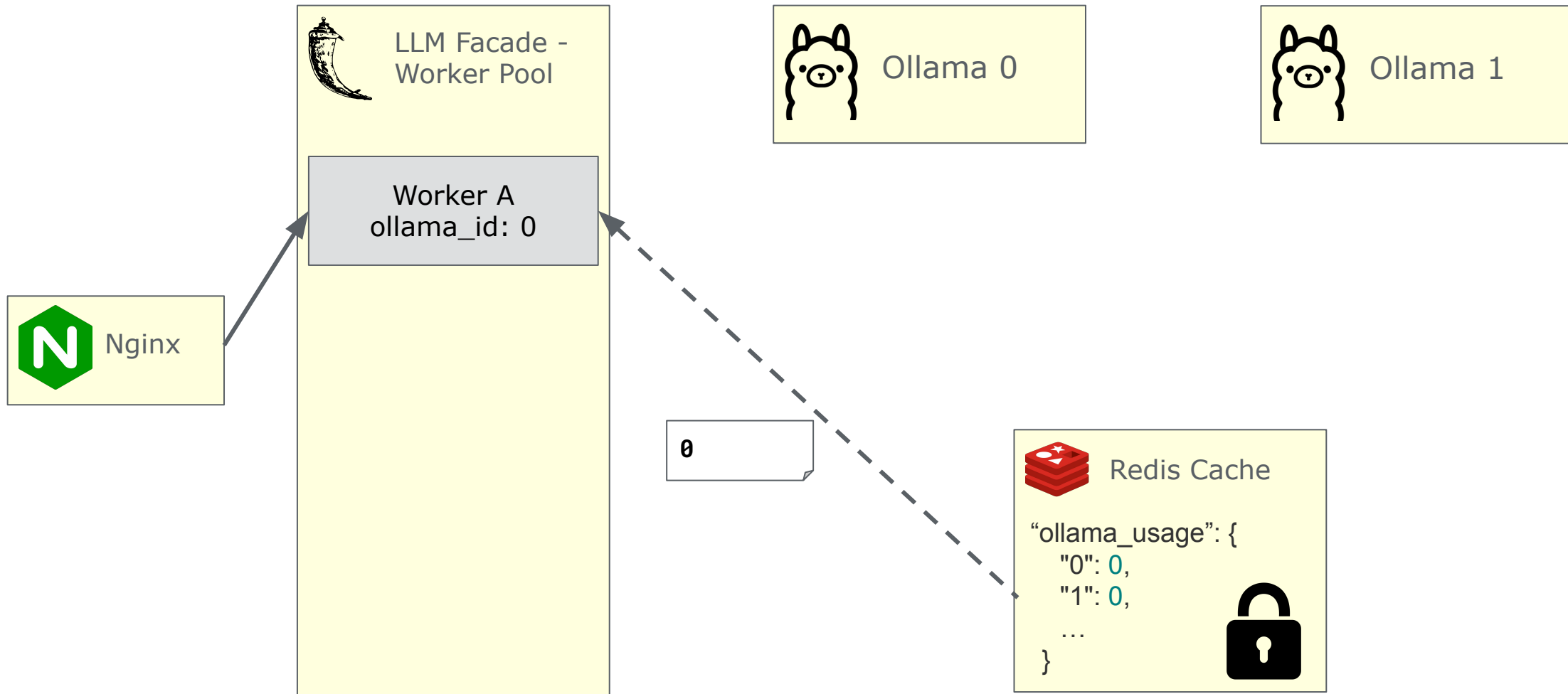
Simultaneous Example Generation



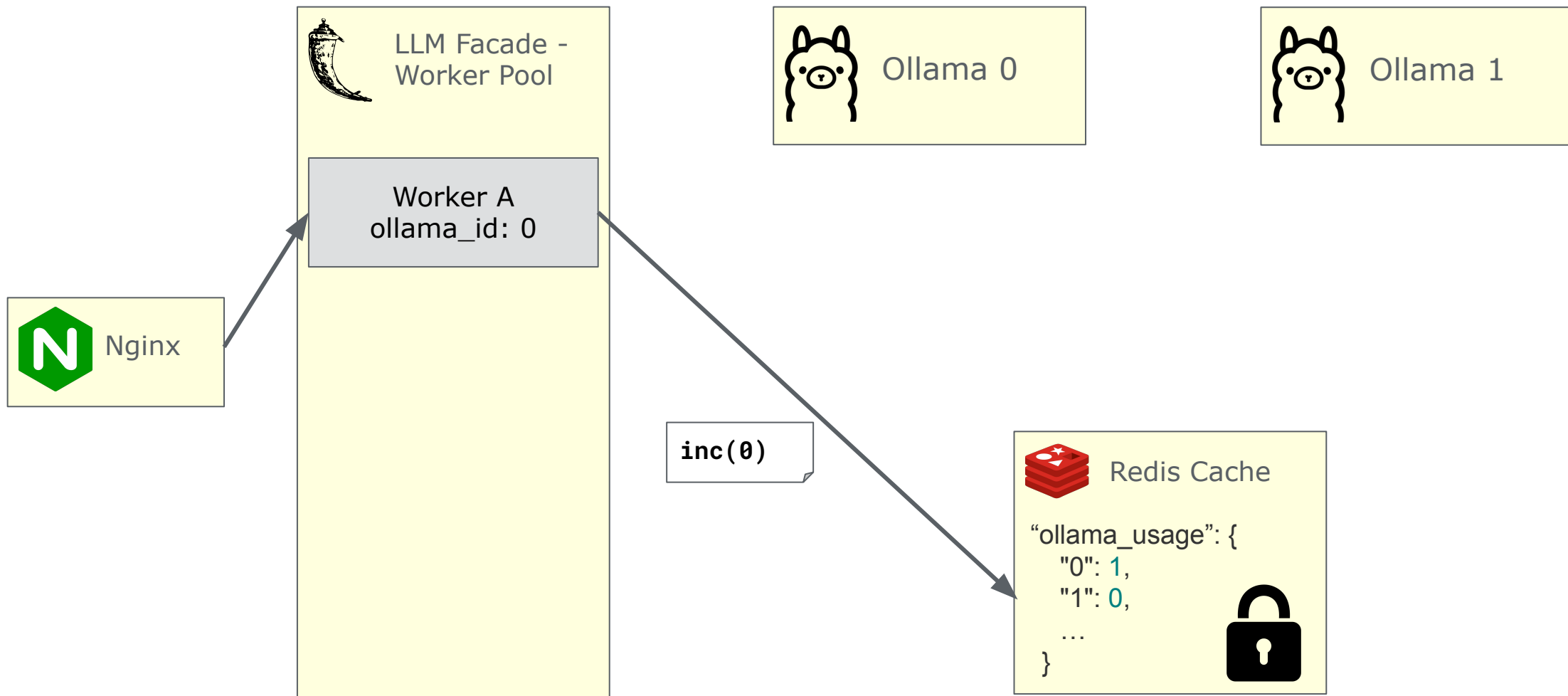
Simultaneous Example Generation



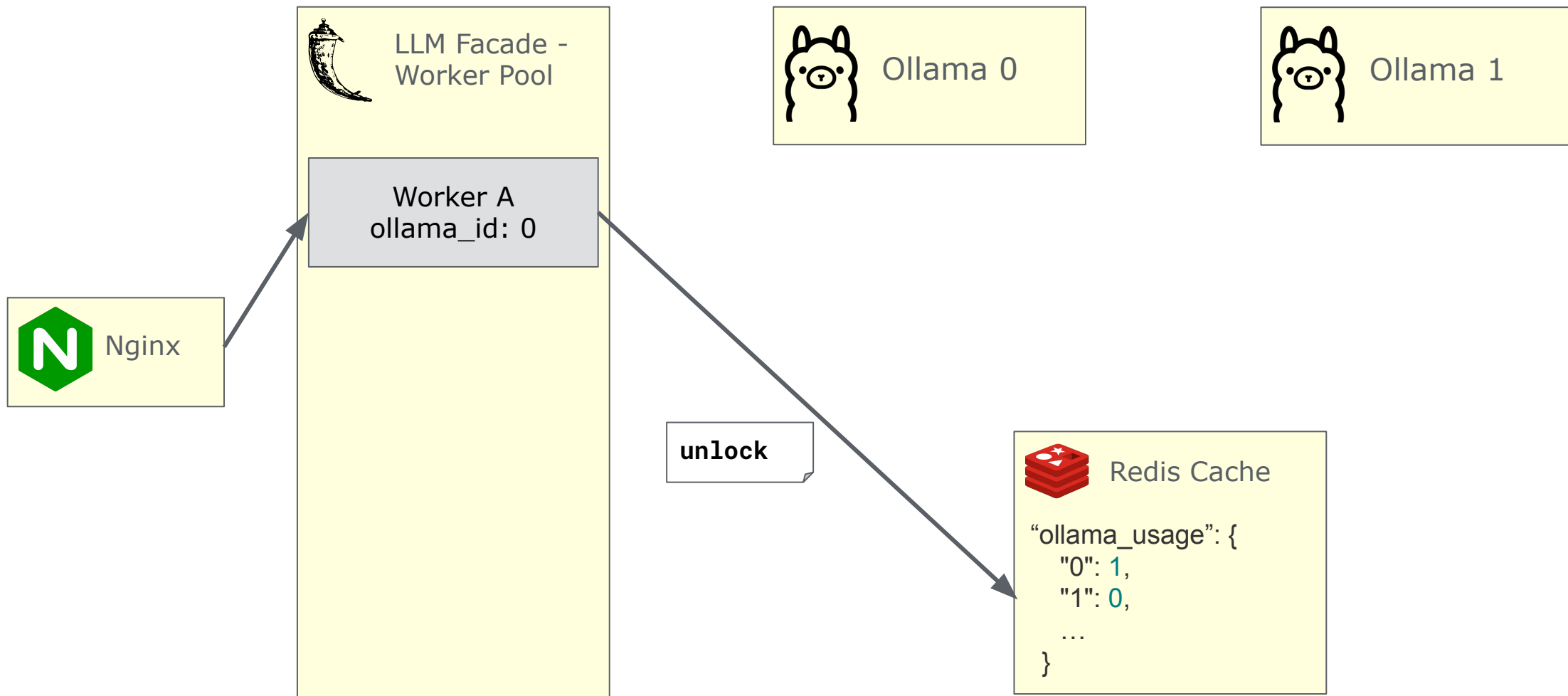
Simultaneous Example Generation



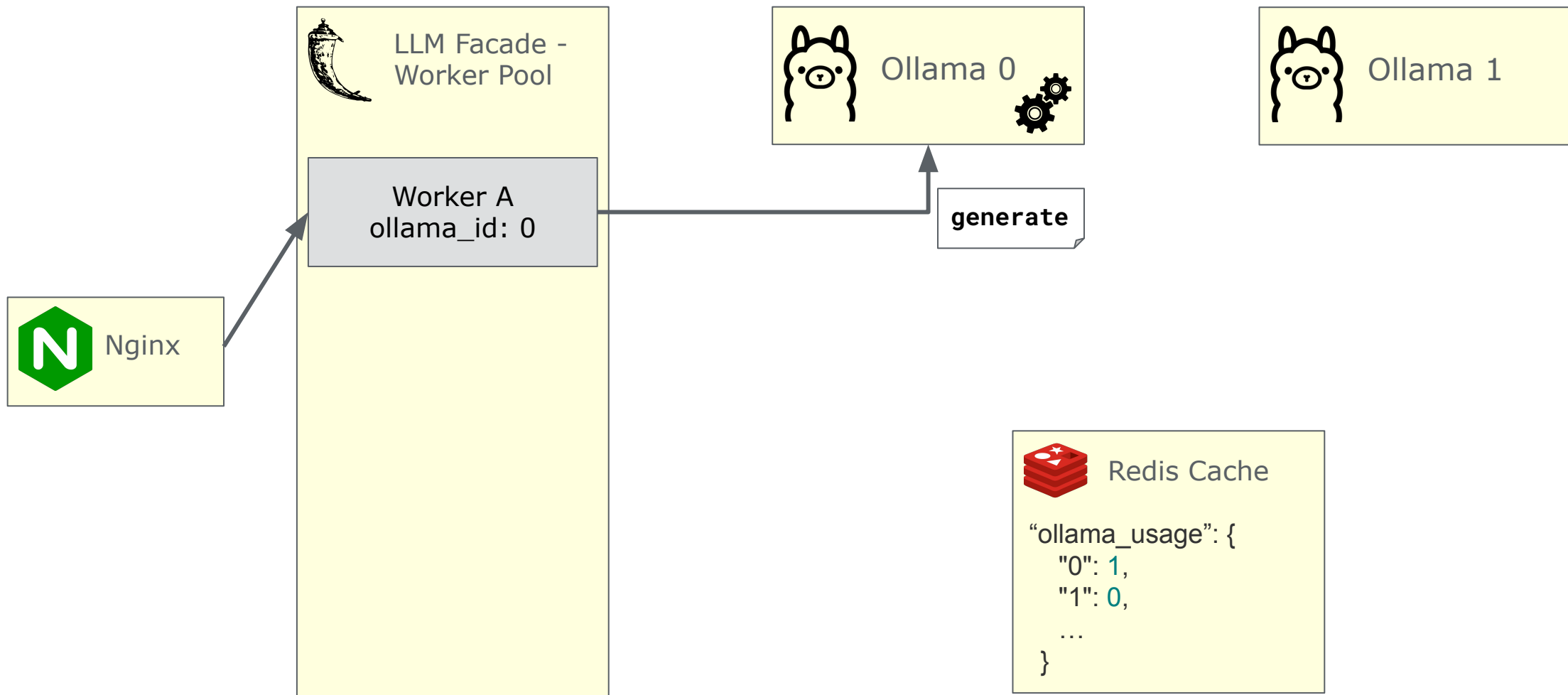
Simultaneous Example Generation



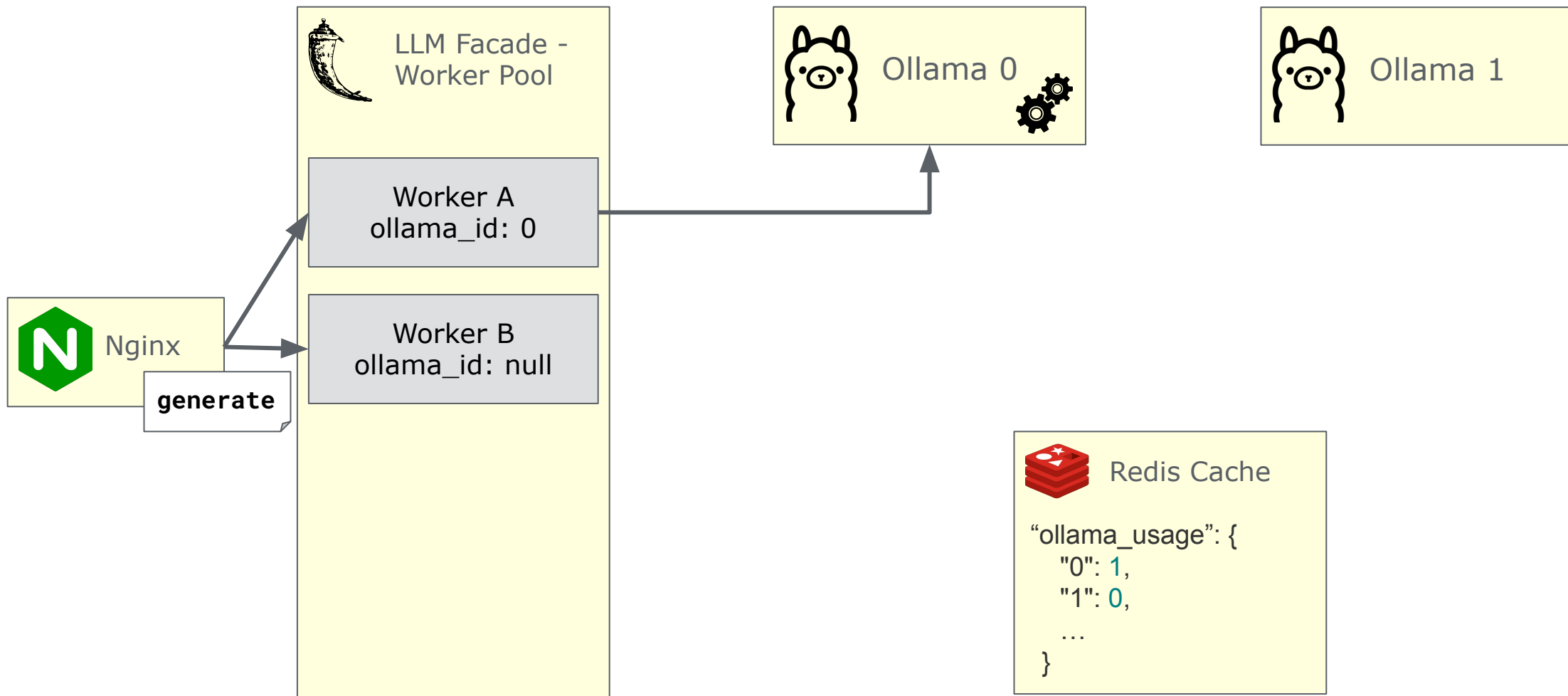
Simultaneous Example Generation



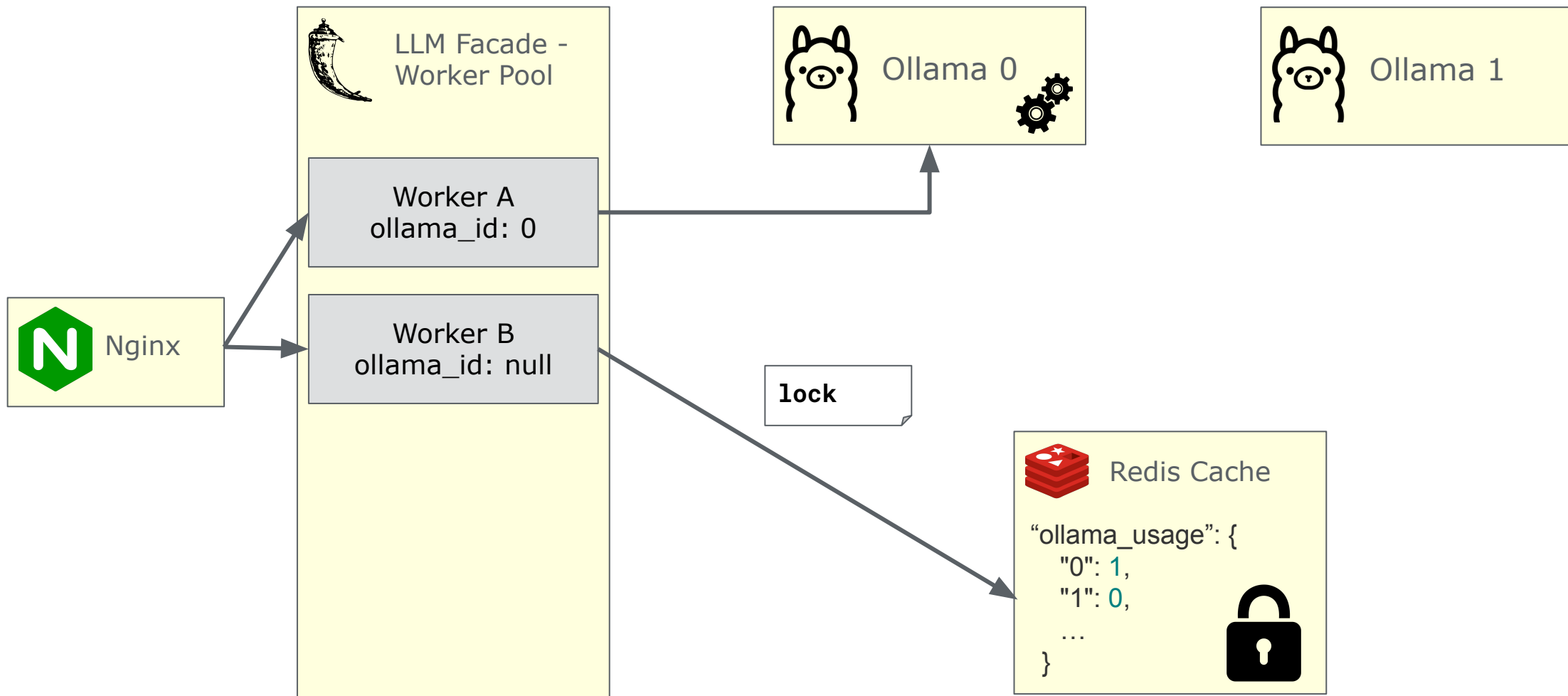
Simultaneous Example Generation



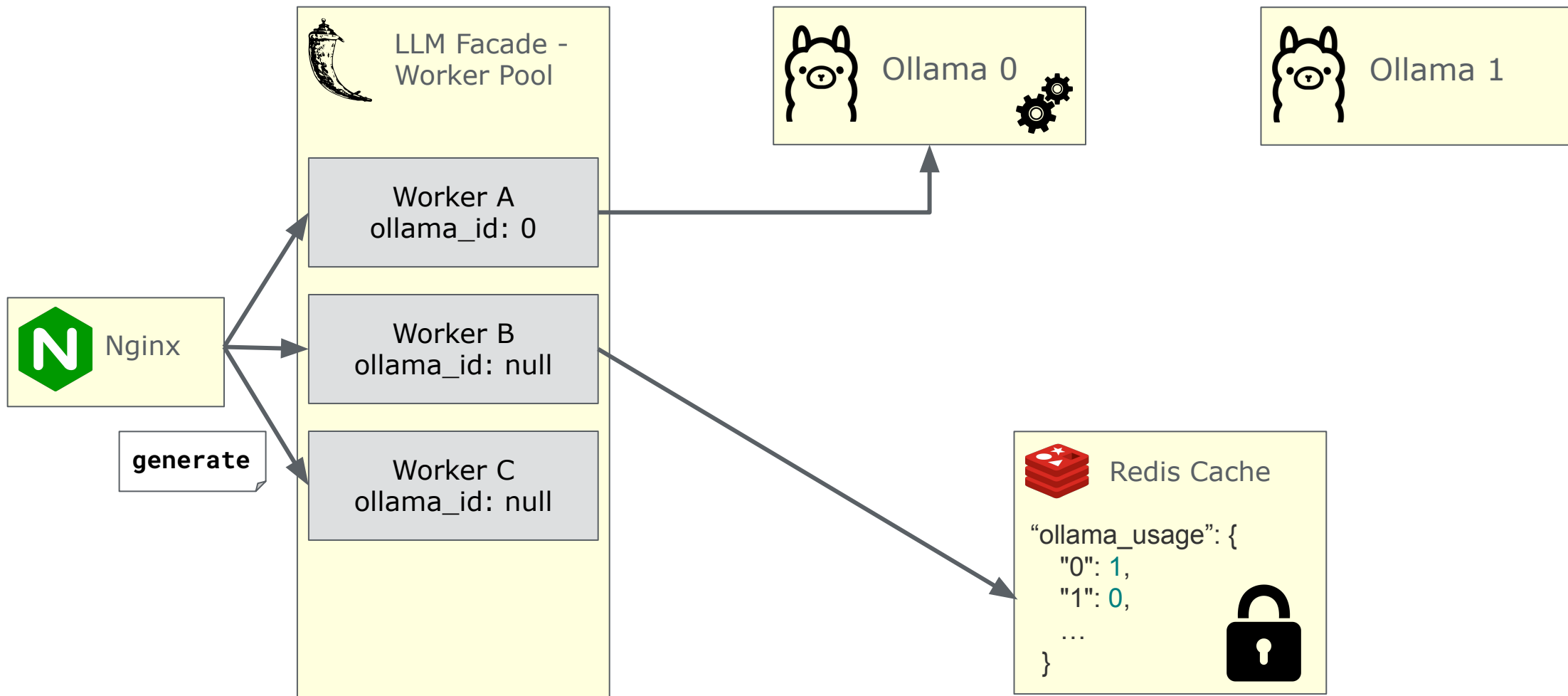
Simultaneous Example Generation



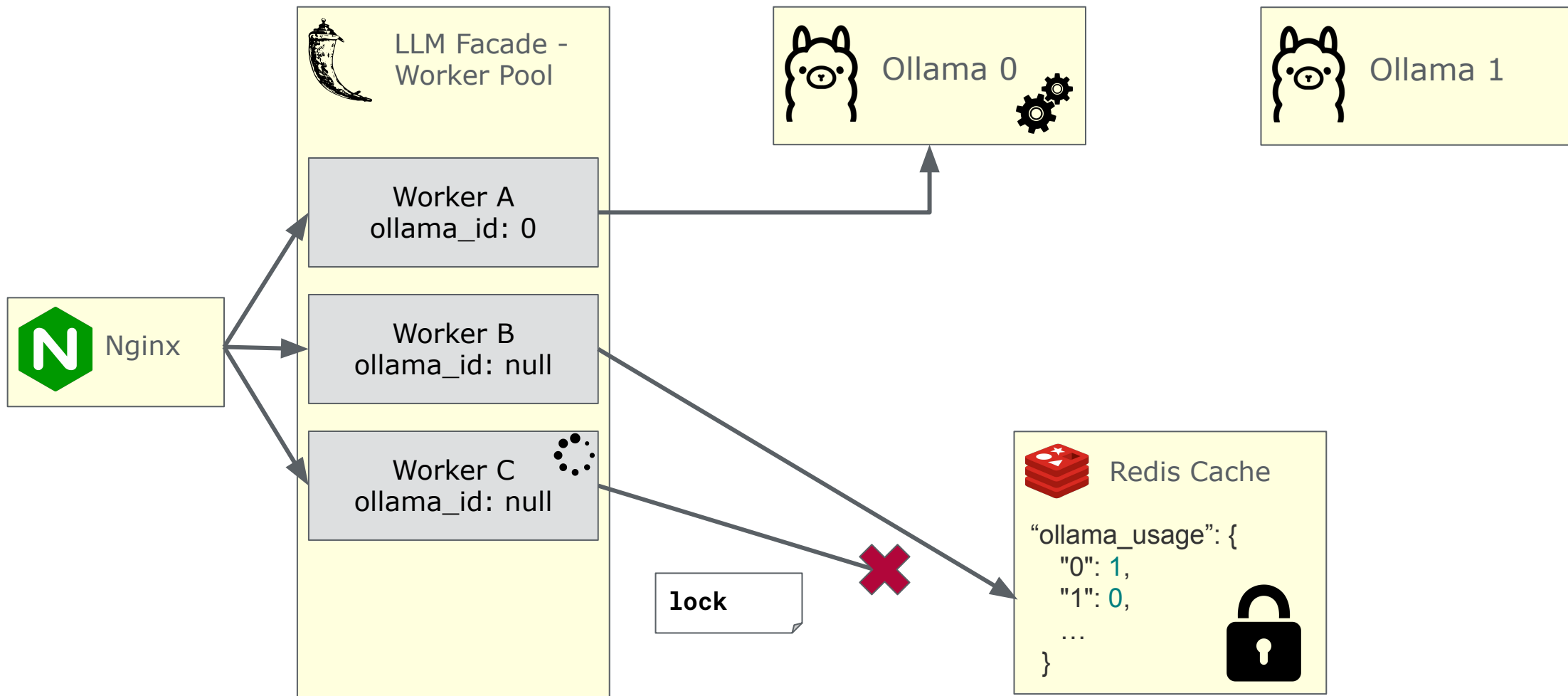
Simultaneous Example Generation



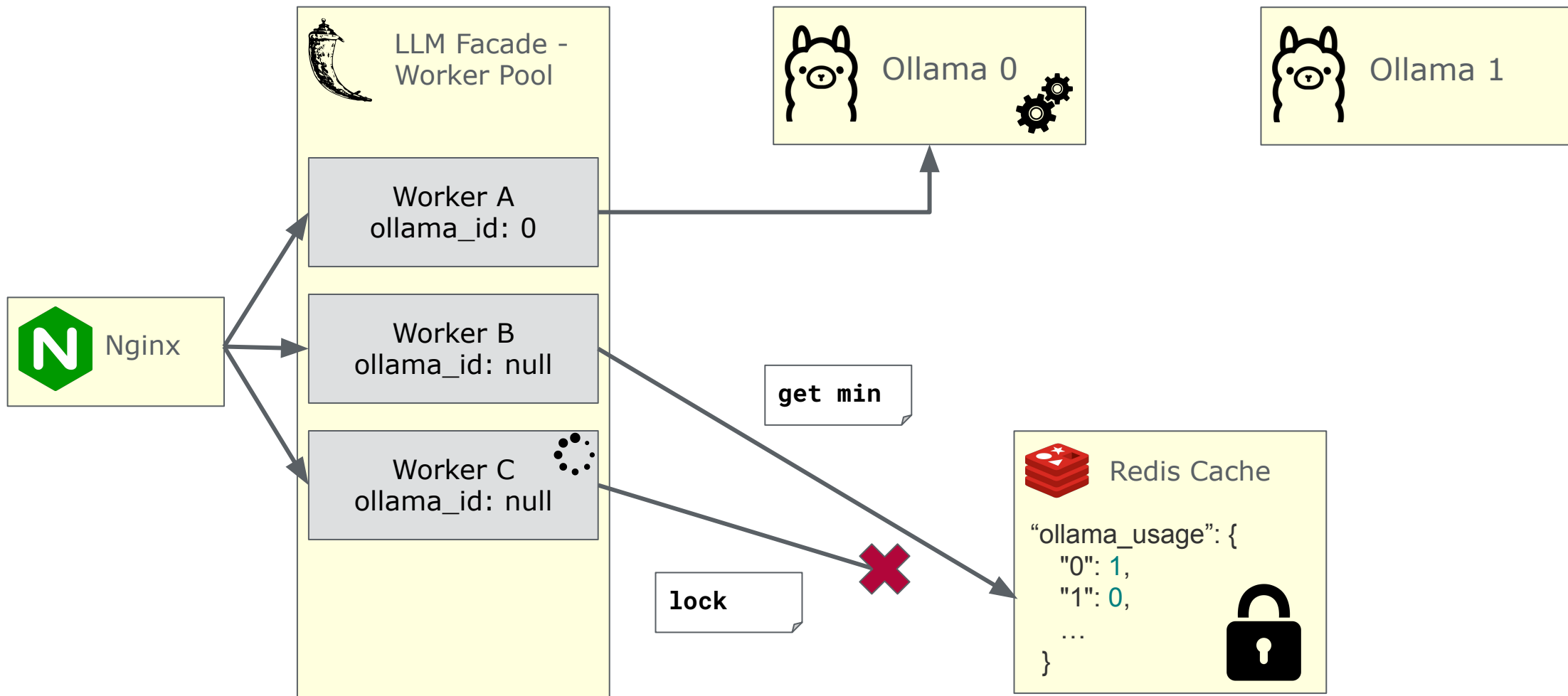
Simultaneous Example Generation



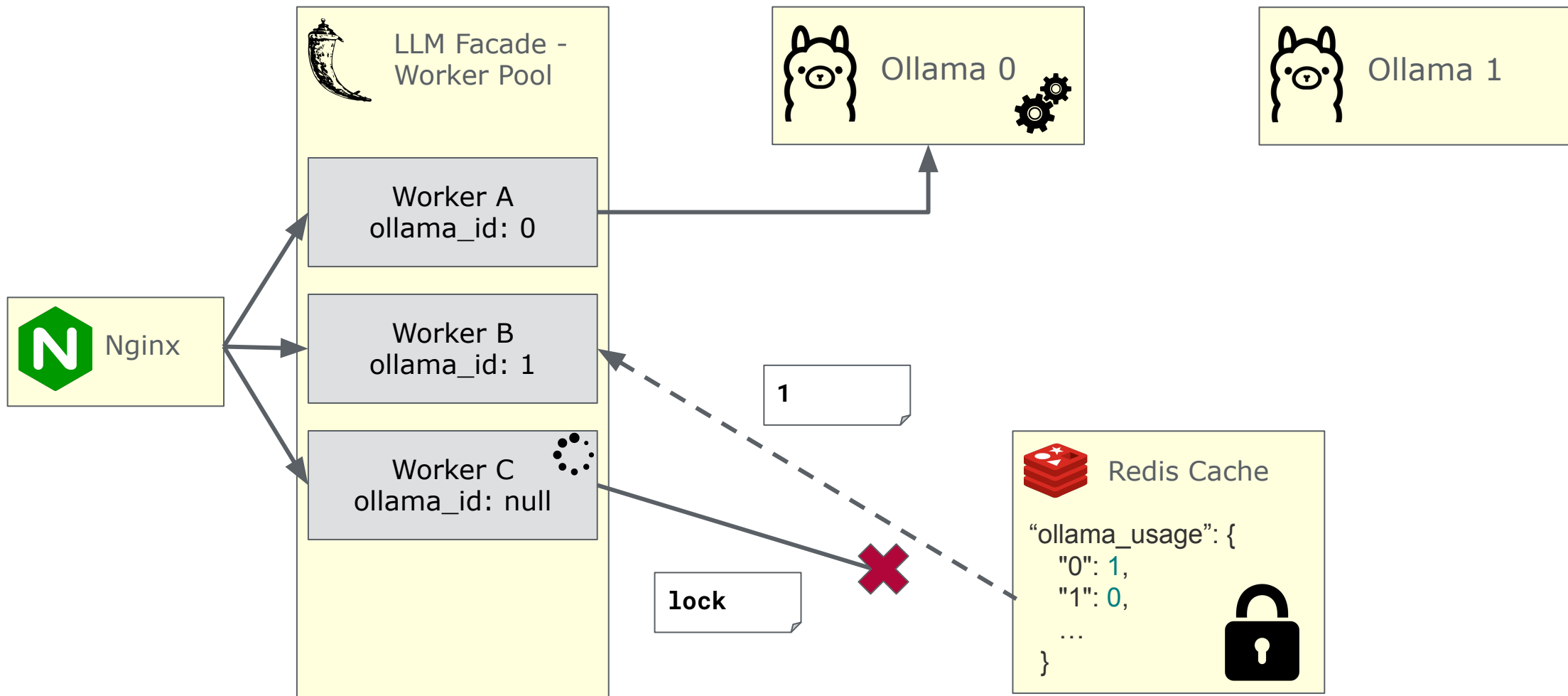
Simultaneous Example Generation



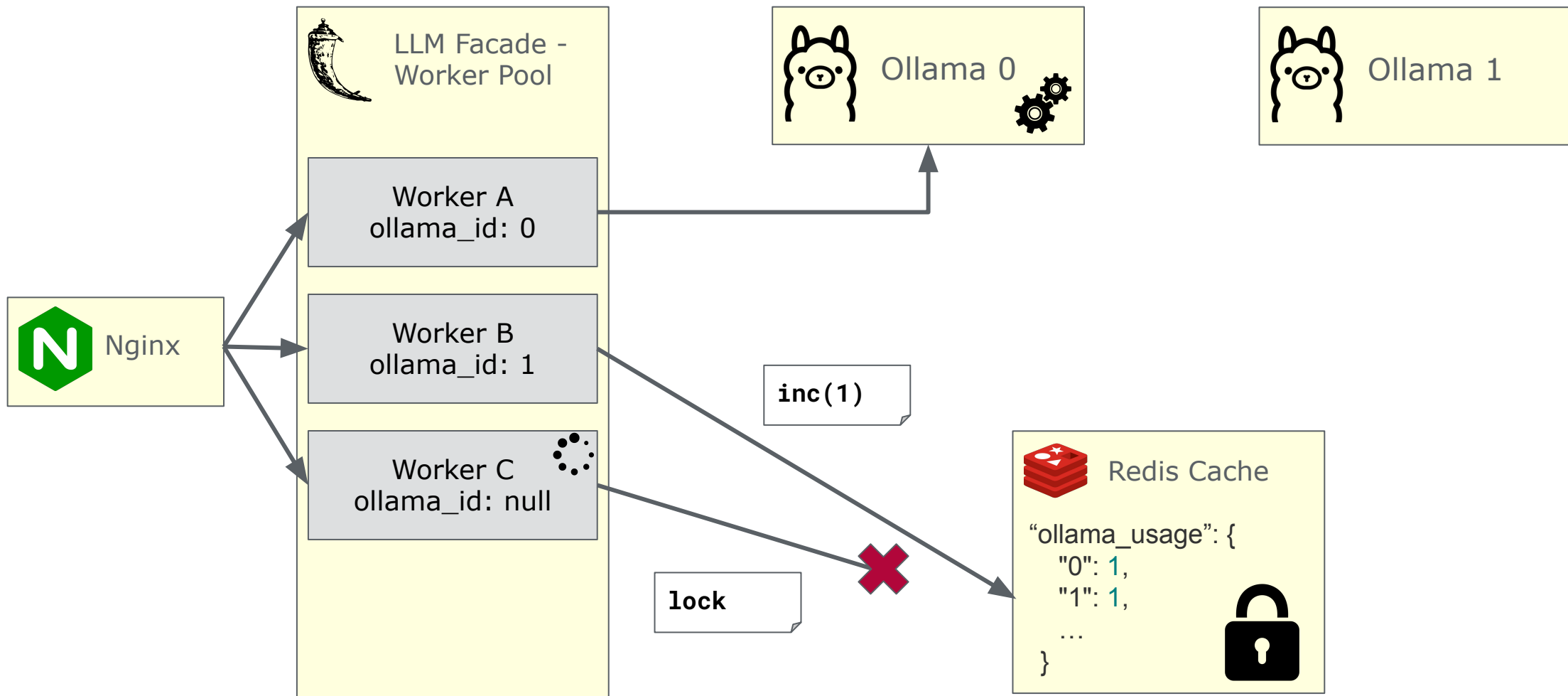
Simultaneous Example Generation



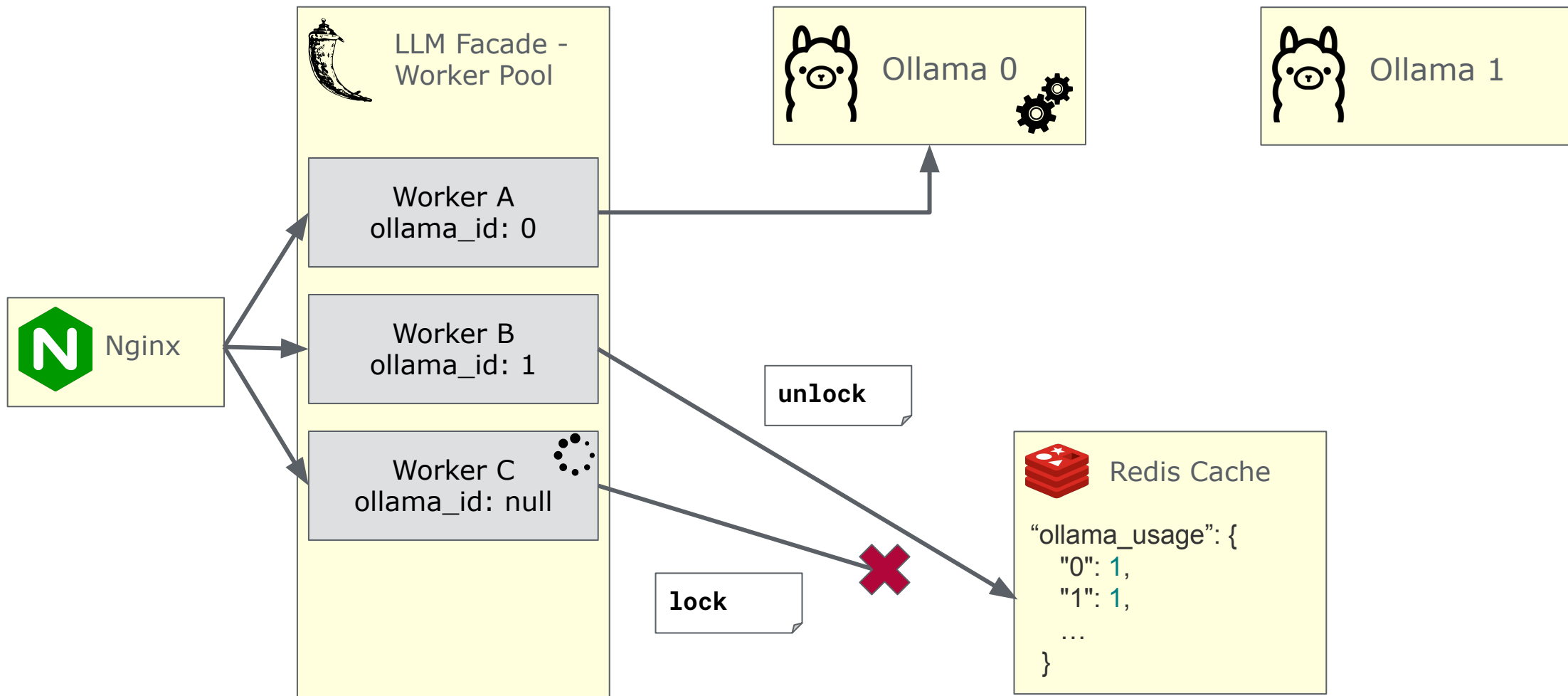
Simultaneous Example Generation



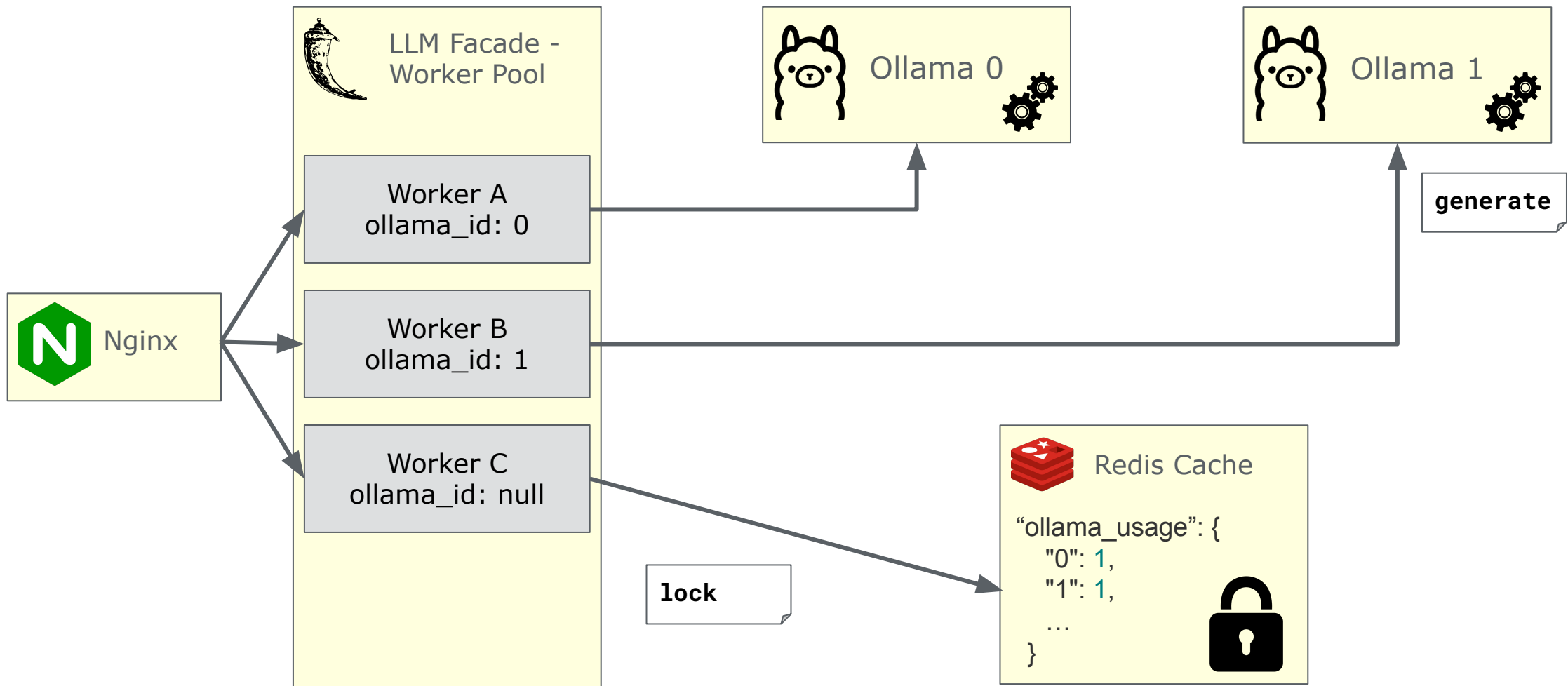
Simultaneous Example Generation



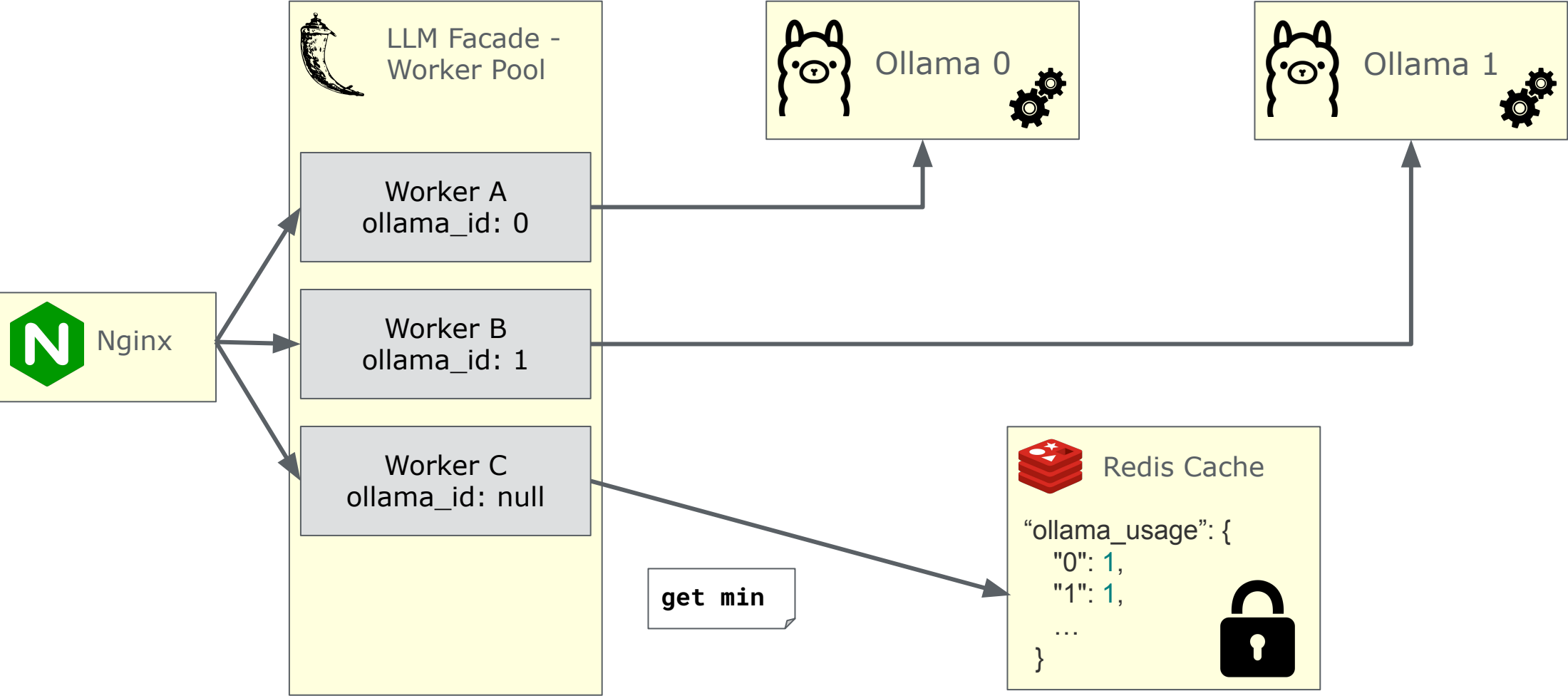
Simultaneous Example Generation



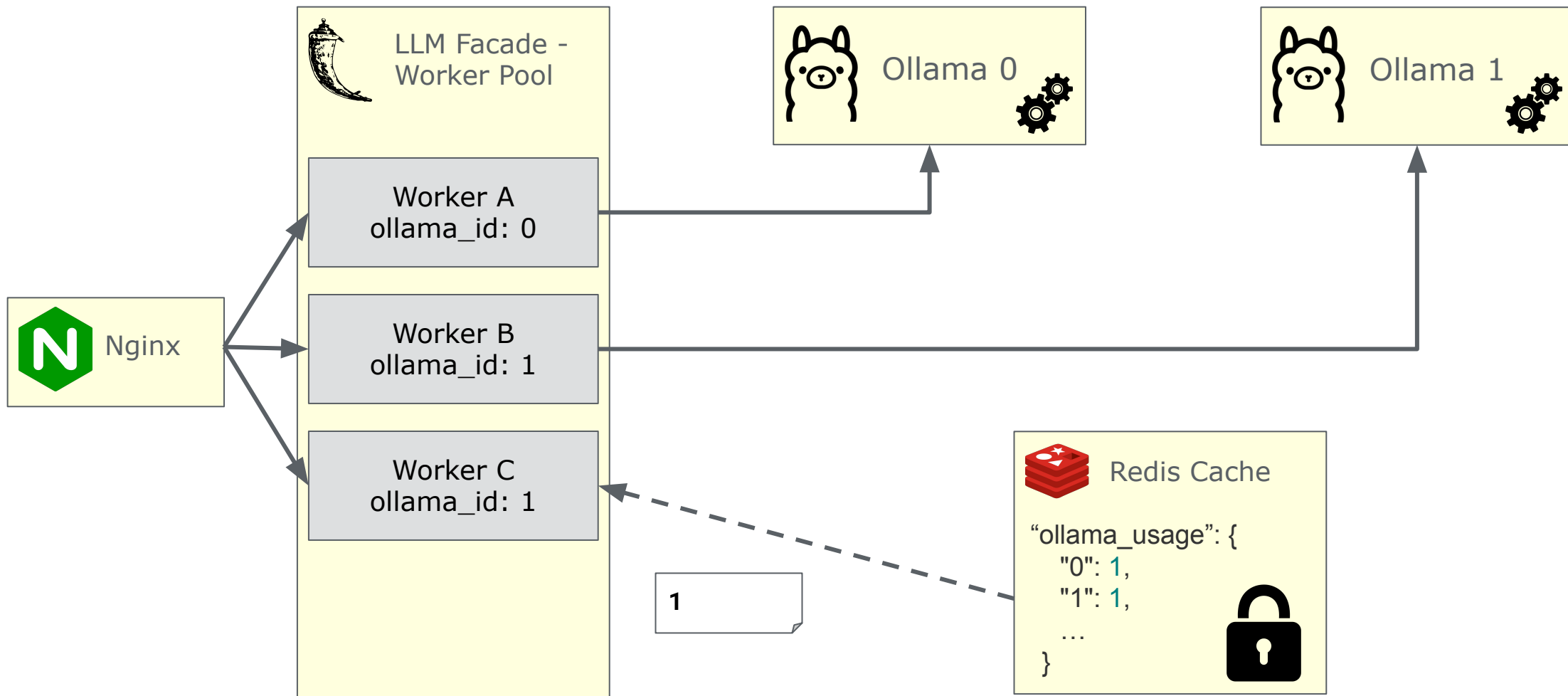
Simultaneous Example Generation



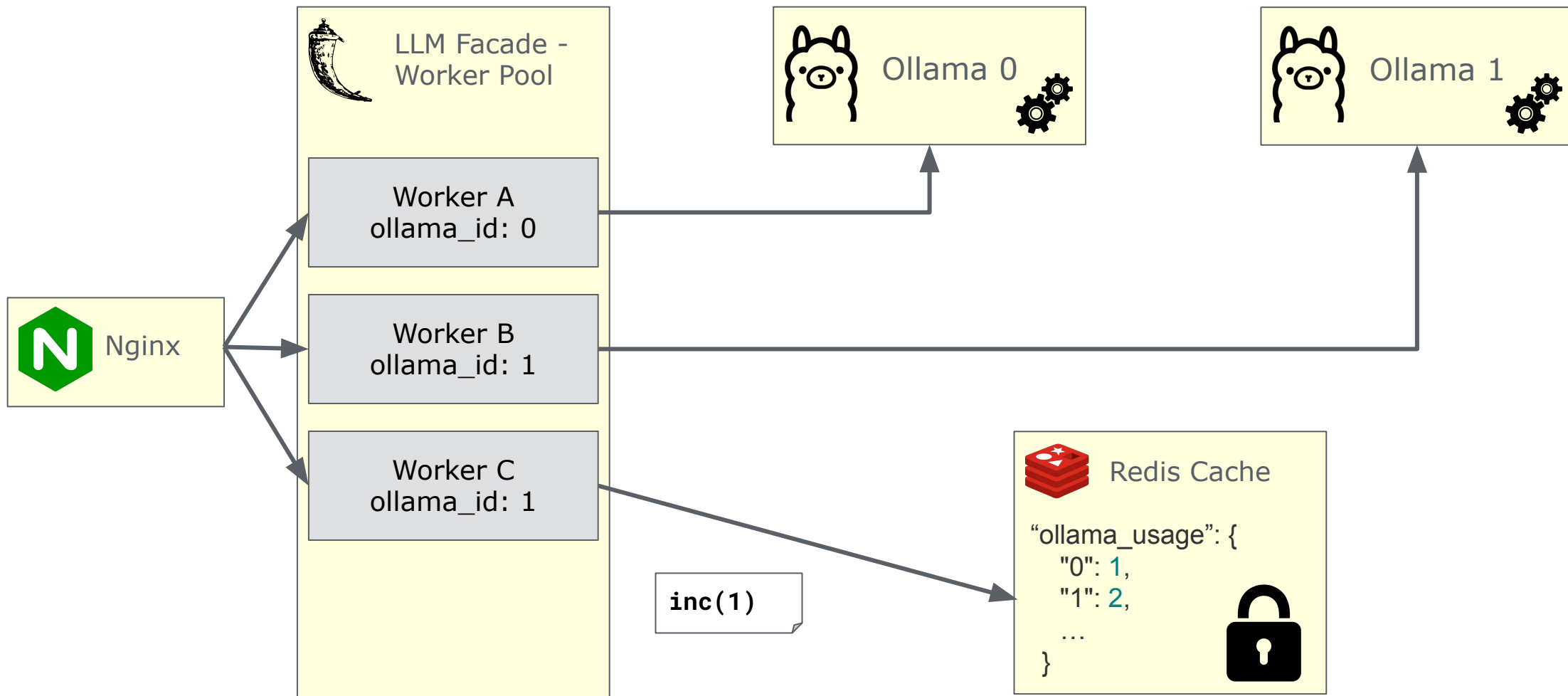
Simultaneous Example Generation



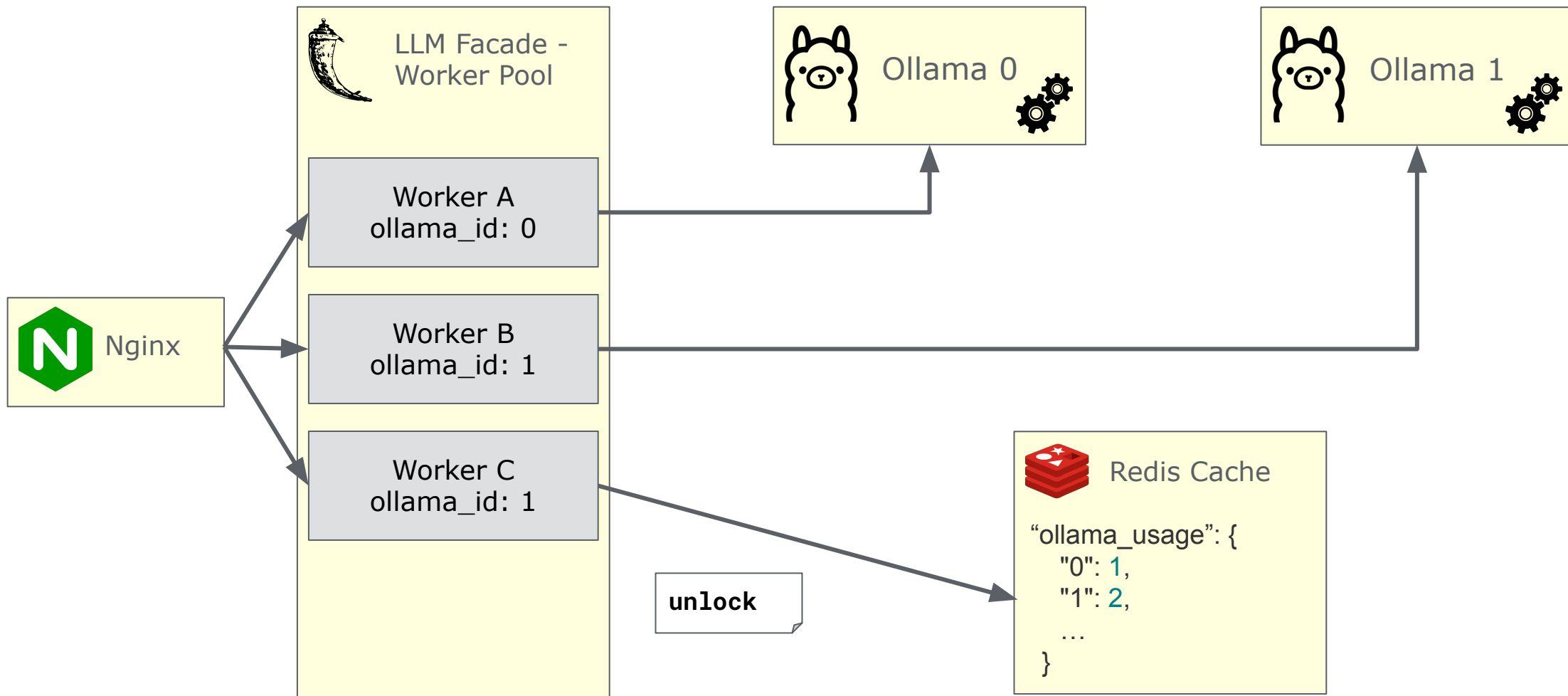
Simultaneous Example Generation



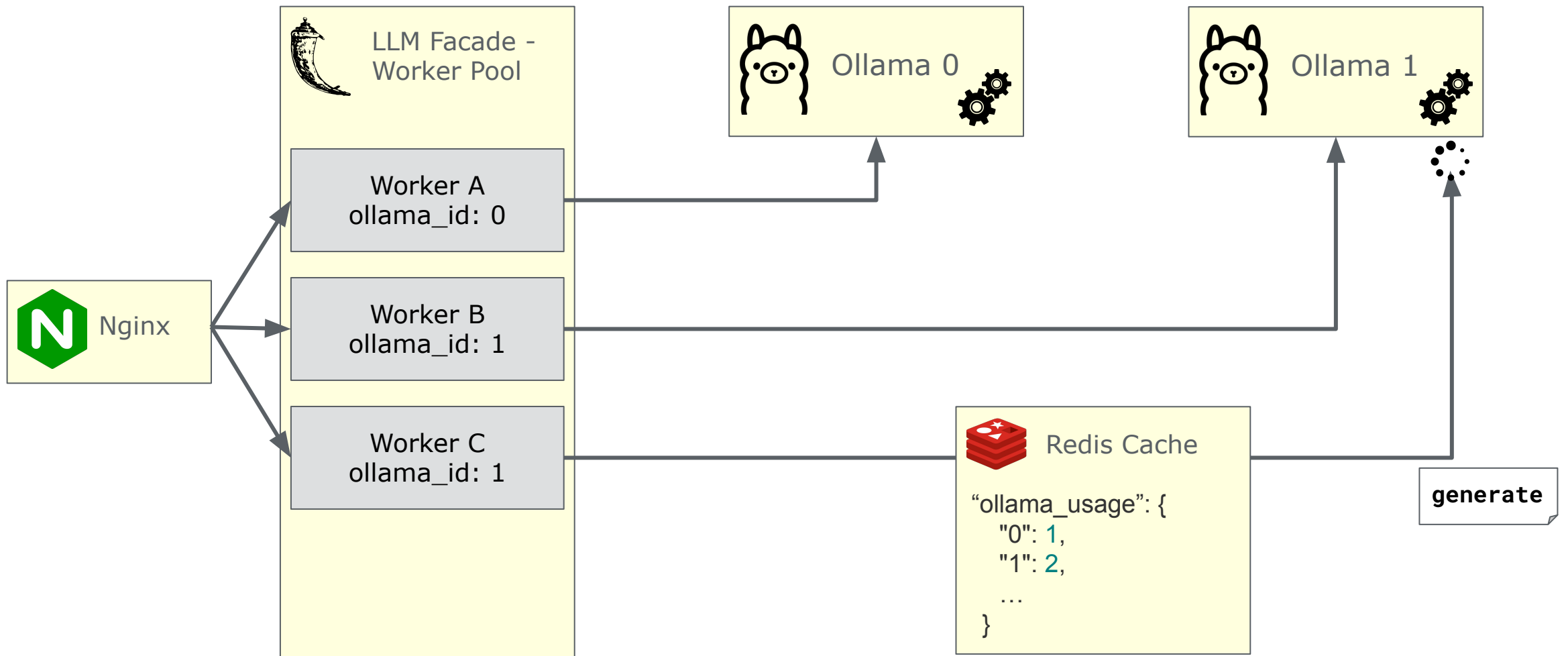
Simultaneous Example Generation



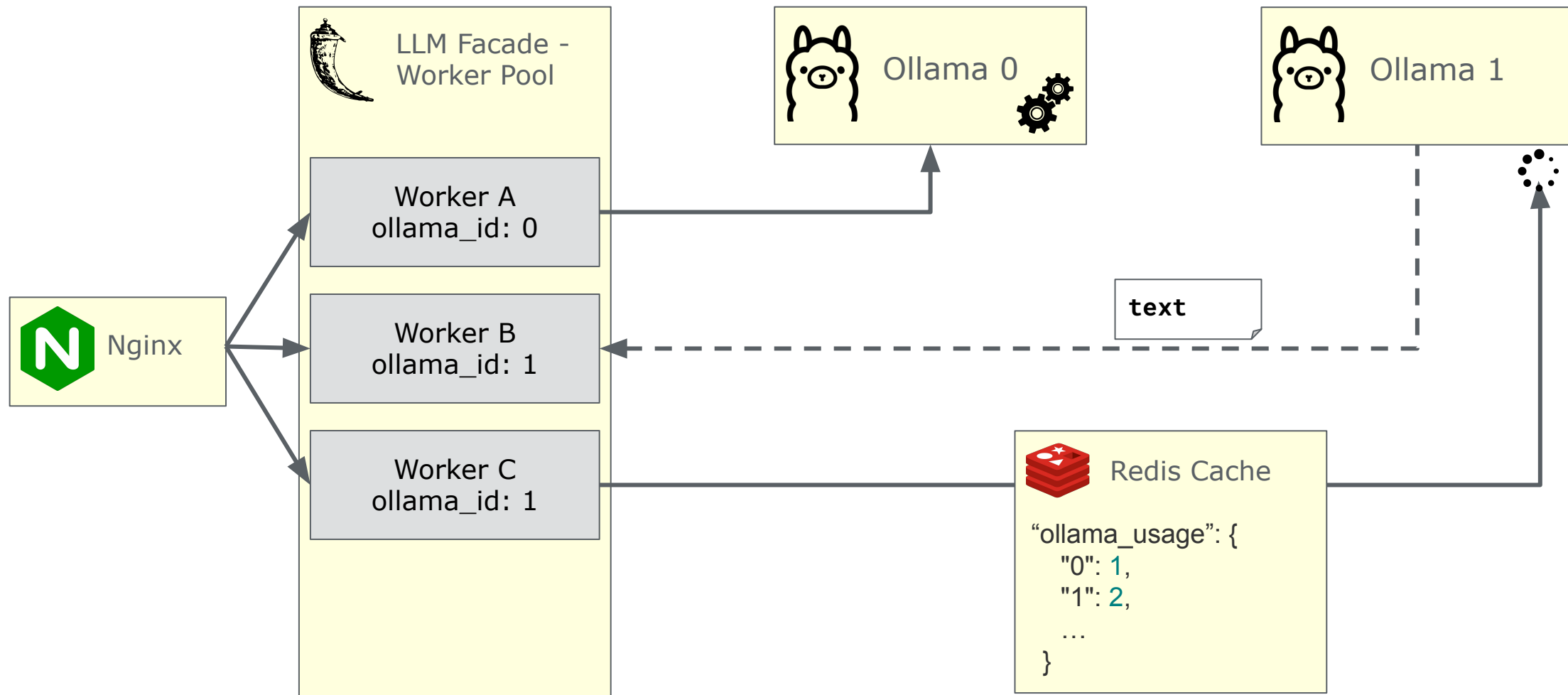
Simultaneous Example Generation



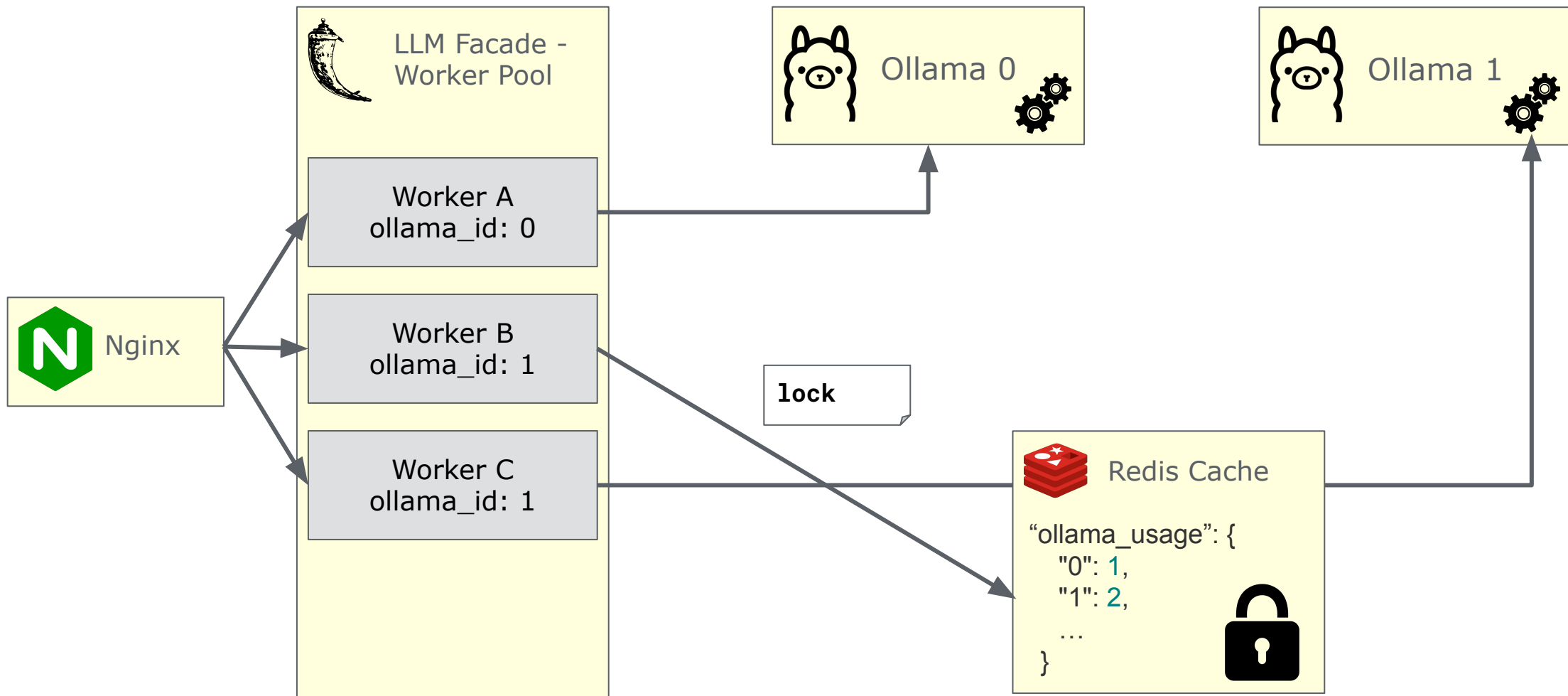
Simultaneous Example Generation



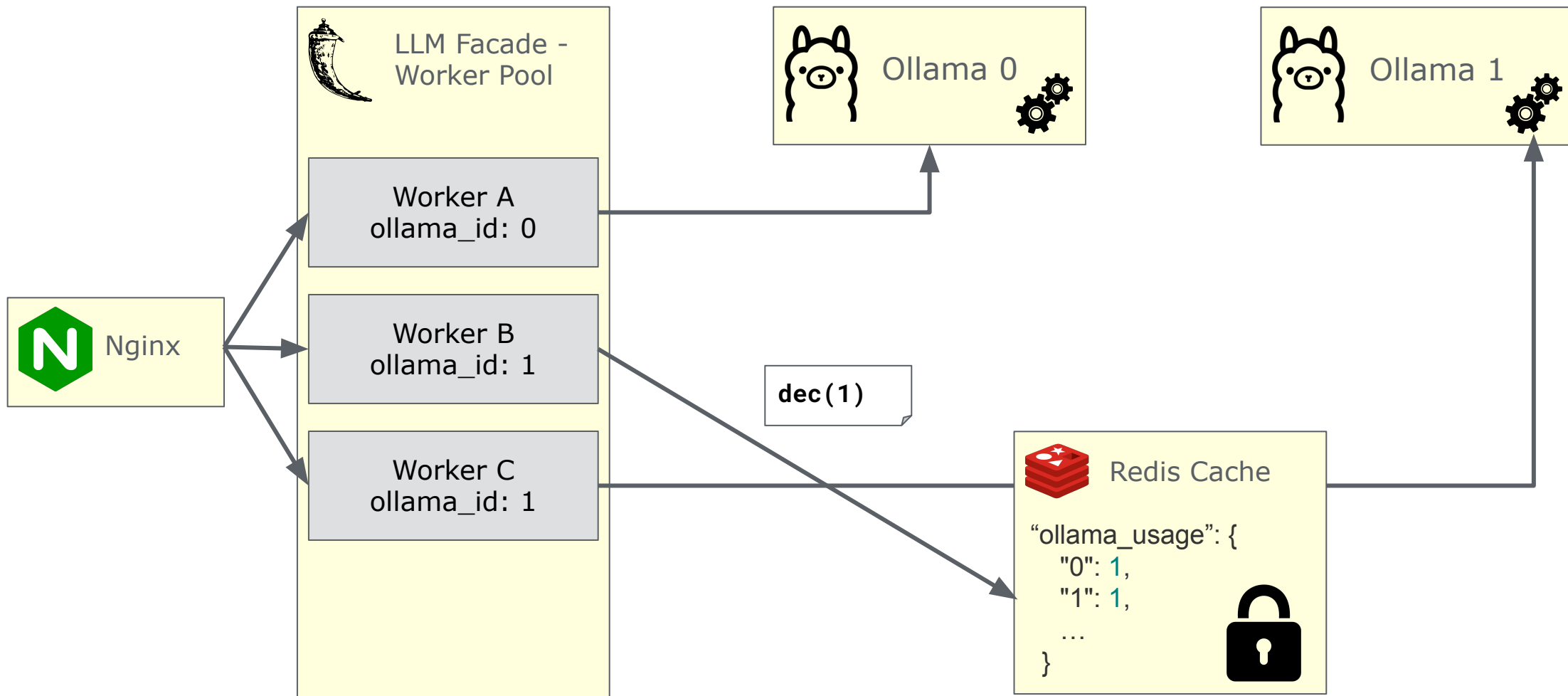
Simultaneous Example Generation



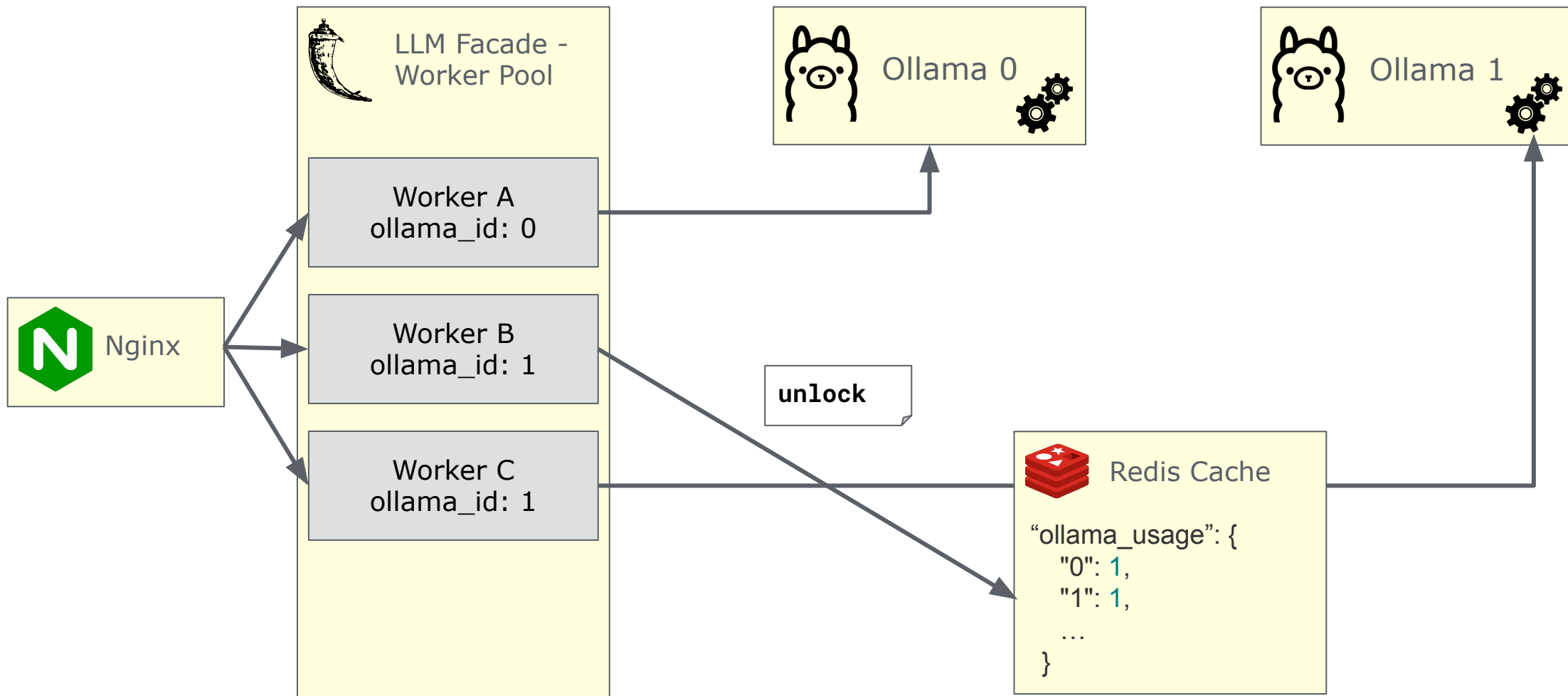
Simultaneous Example Generation



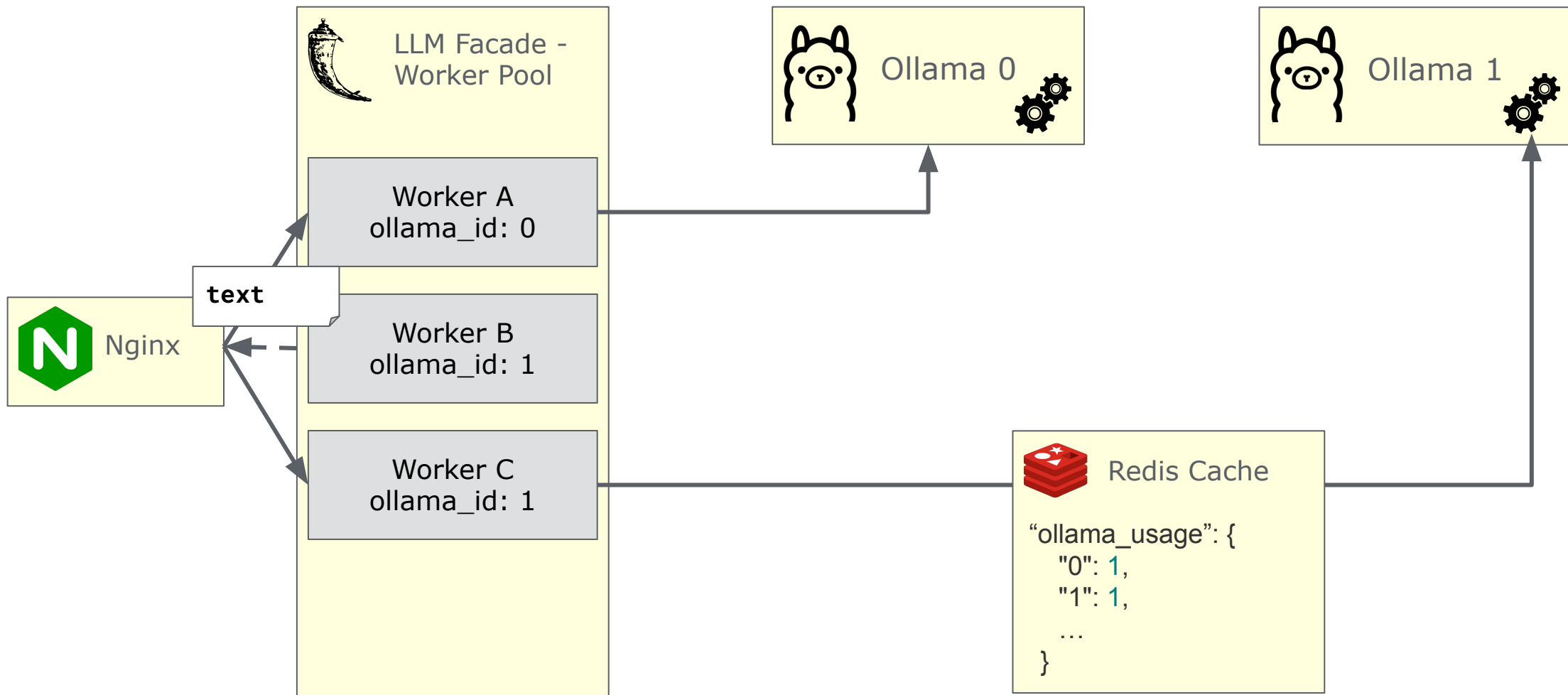
Simultaneous Example Generation



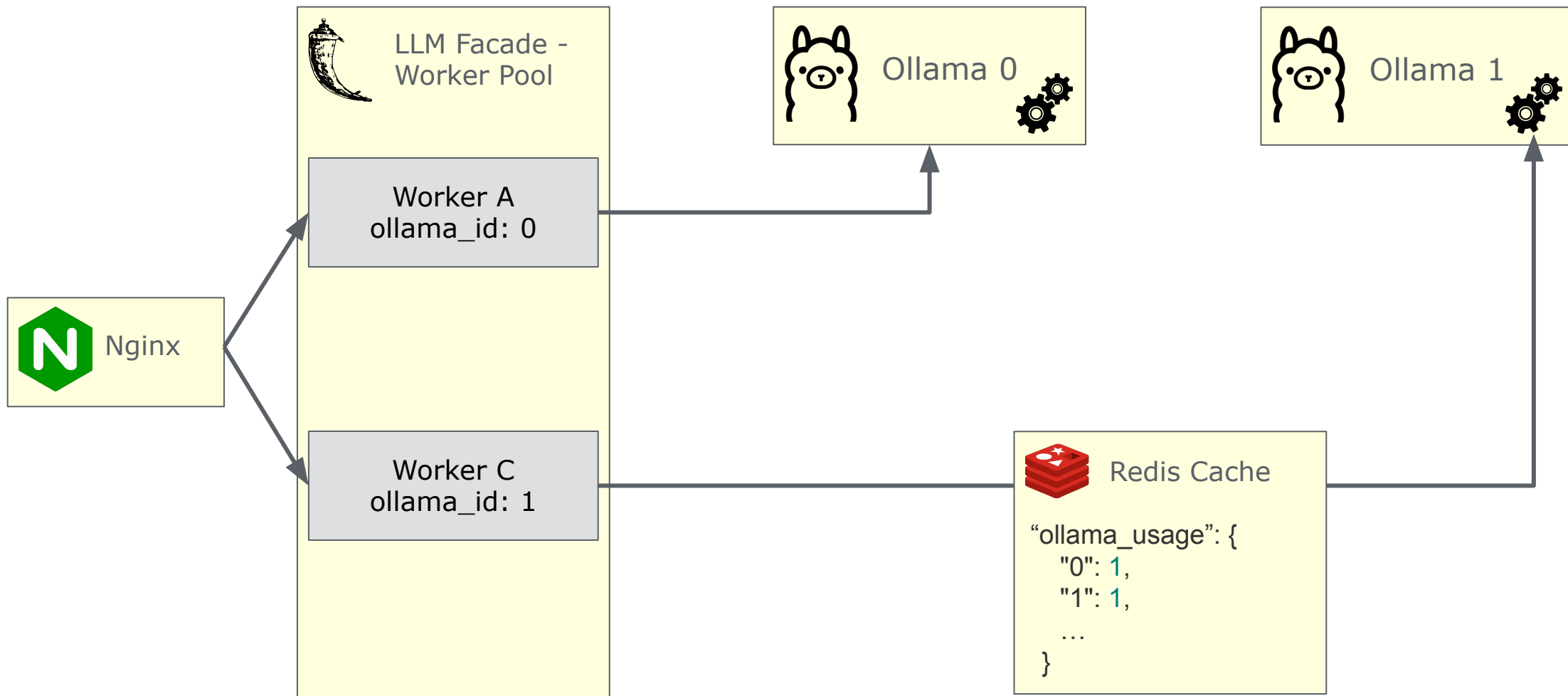
Simultaneous Example Generation



Simultaneous Example Generation



Simultaneous Example Generation



Challenges

Code Related

- Simultaneous Generation
- Syntax Highlighting
- Prompt Engineering to avoid explanation in code
- Hard to find suitable metrics
 - Code Length
 - Code Complexity
 - Maintainability

Project Related

- User Feedback: Server restricted to HPI network
- Linux and Windows incompatibilities

Challenges

Code Related

"You are as experienced as a senior software developer. You will get a piece of code documentation and the source code for this documentation.

Your task is to provide a helpful code example for the documentation. You are encouraged to use also other functions from the same framework to demonstrate interoperability between them if it is appropriate.

You ONLY output python code. If you output anything else than python code 100 cute little rabbits will DIE. You DON'T under any circumstances output any explanation of the code."

- Maintainability

Threatening the model helped reduce the output of any text other than code, but could not eliminate it entirely.

Challenges

Code Related

- Simultaneous Generation
- Syntax Highlighting
- Prompt Engineering to avoid explanation in code
- Hard to find suitable metrics
 - Code Length
 - Code Complexity
 - Maintainability

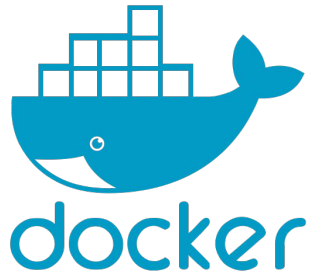
The only metric, which is actually meaningful, is a user's opinion of the example and whether it's helpful to them or not.

Project Related

- User Feedback: Server restricted to HPI network
- Linux and Windows incompatibilities

Learnings / Reflection

New Technologies



ollama



Project Related

- Improve Collective Code ownership
- Two development phases
- Focus on user experience / story

User Feedback

Methodology

- Introduction to subject
- ~15 Minutes
- Communicated incomplete features

We asked several users for feedback in the last 1-2 weeks of development. After giving them a brief introduction to the subject, i.e. telling them that we built an application to generate code examples for different APIs, we let them explore the tool on their own.

Two constraints made the collection of user feedback slightly challenging: The website was only available from the HPI network and we were not able to find users which were currently in a learning process for any of the available APIs. Therefore, we were limited to brief, in-person sessions of about 15 minutes, gathering mostly “theoretical” feedback. The users were asked to generate a few examples, experiment with the workflow and use as many features as possible. When they were using any incomplete features, we made them aware of this.

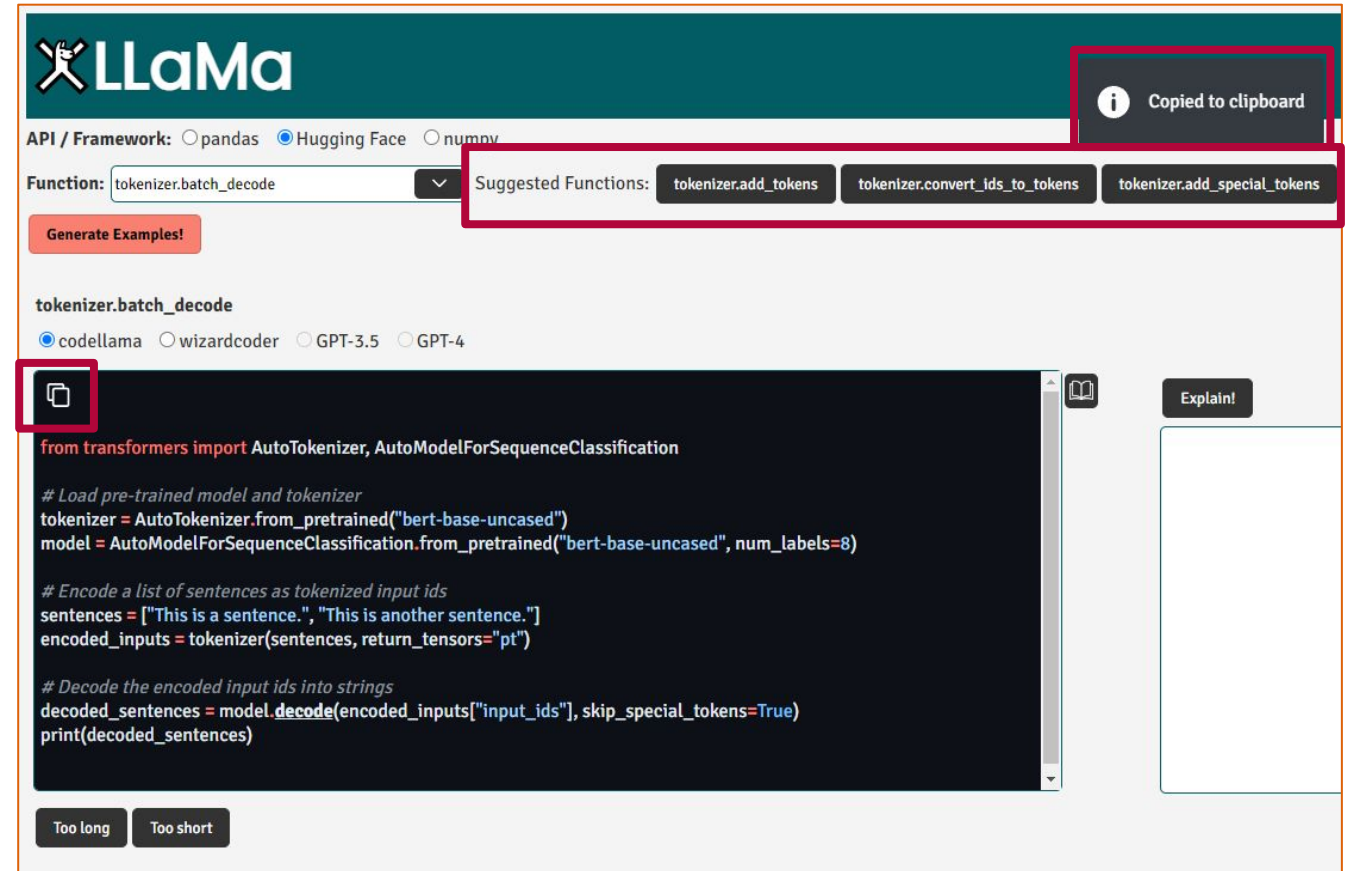
User Feedback

Methodology

- Introduction to subject
- ~15 Minutes
- Communicated incomplete features

Implemented

- Confusing function buttons
- Copy button should give visual feedback



User Feedback

User 1 uses pandas for data analysis.

User 1 (Molecular Biology PhD Student)

- Documentation button very helpful
- Did not notice clickable button in example
- More advanced options for generated examples and layout
- History of generated examples
- Would even just use it to read documentation, nice collection of different frameworks

User 2 (Machine Learning Engineer)

- Currently learning new framework with bad documentation, so tool would be useful if all functions were included

User 3 (Software Developer)

- Would use when learning new API
- Would be nice to have for programming languages, e.g. Java
- Nice to have everything in one place

Idea

Chatbots

- Helpful when working and learning new tools, but often cumbersome to use
 - Might need long prompts ✓
 - First results may be unhelpful (✓)
 - Multiple follow up prompts can be annoying ✓
- Different LLMs give different results (of different quality) ✓

Problems like long prompts or multiple follow up prompts have been eliminated, since xLLaMa uses no prompts whatsoever. Since examples are generated with multiple LLMs at the same time, the change of having a helpful result immediately are higher.



ChatGPT

The `get_added_vocab` method in Huggingface's Transformers library is used to retrieve the tokens added during the tokenizer's training from a specific

Here's an example of h

You

And how would i now use `prepare_for_tokenization`?

ChatGPT

The `prepare_for_tokenization` method in Huggingface's Transformers library allows you to preprocess text before tokenization. This method is useful for ensuring

You

I don't understand the third line

ChatGPT

My apologies for the confusion. Let me clarify the third line of code.

Users are now able to generate examples for poorly documented APIs, which might not provide examples on their own. However, they are also able to find additional examples for well documented APIs, in a fast and convenient way.

APIs

- Some API references are very detailed and helpful, others are not ✓

Future Work

- Add entire APIs to database
- Add Context integration
- Add GPT or any other LLM support
- Collect and use user feedback

It would be nice to collect and store user feedback, specifically about the helpfulness of an example, store it and use this data to fine-tune a model and generate even more helpful examples in the future.

The screenshot displays the LLaMa API interface. At the top, the 'API / Framework' section has 'Hugging Face' selected. The 'Function' dropdown is set to 'tokenizer.batch_decode'. Below this, there are buttons for 'Generate Examples!', 'tokenizer.convert_ids_to_tokens', 'tokenizer.add_tokens', 'tokenizer.add_special_tokens', and 'tokenizer.apply_chat_template'. The main content area shows the 'tokenizer.batch_decode' function with 'codellama' selected. A code editor displays a Python script for loading a BERT model and tokenizer, encoding sentences, and decoding the input IDs. To the right of the code editor is an 'Explain!' button and a text box containing a step-by-step explanation of the code. At the bottom of the code editor are 'Too long' and 'Too short' buttons.

```
from transformers import AutoTokenizer, AutoModelForSequenceClassification

# Load pre-trained model and tokenizer
tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=8)

# Encode a list of sentences as tokenized input ids
sentences = ["This is a sentence.", "This is another sentence."]
encoded_inputs = tokenizer(sentences, return_tensors="pt")

# Decode the encoded input ids into strings
decoded_sentences = model.decode(encoded_inputs["input_ids"], skip_special_tokens=True)
print(decoded_sentences)
```

1. from transformers import AutoTokenizer, AutoModelForSequenceClassification
from the Hugging Face Transformers library

2. tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
assigns it to the variable tokenizer

LLaMa

Thank you for your attention!

Links:



<http://delos.eaalab.hpi.uni-potsdam.de:8080/>



<https://github.com/CR1337/xLLaMa>



Quellen

- <https://pandas.pydata.org/docs/>
- <https://numpy.org/doc/>
- <https://huggingface.co/docs>
- <https://ollama.ai/>
- <https://huggingface.co/codellama/CodeLlama-70b-hf>
- <https://huggingface.co/WizardLM/WizardCoder-33B-V1.1>
- <https://chat.openai.com/>