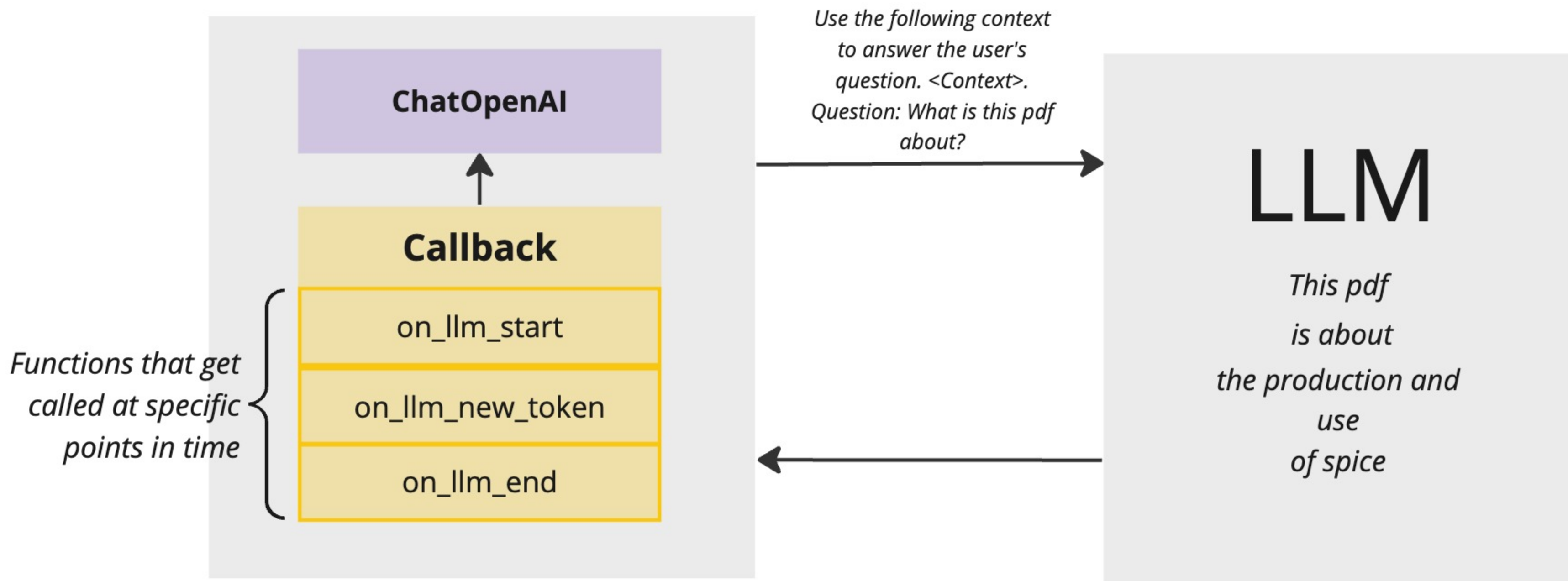
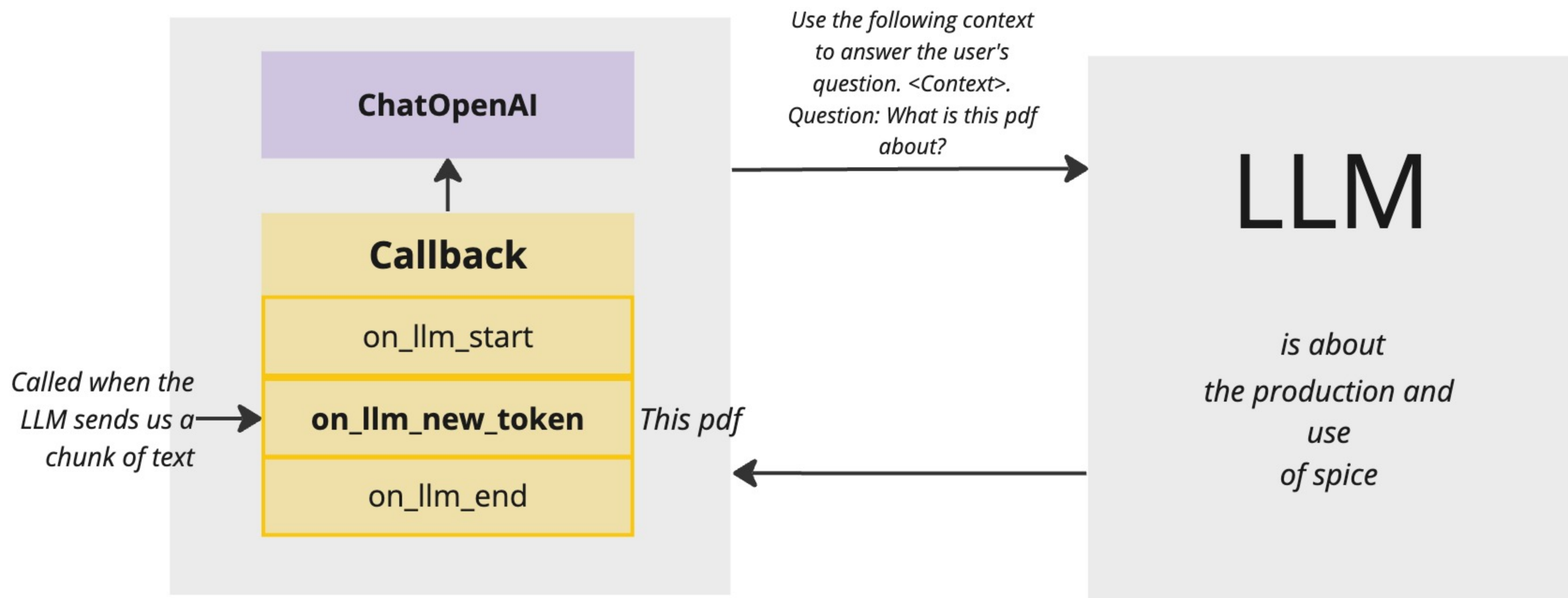


**Our app will support both sync + streaming text generation (just for demo purposes)**



Need to figure out how to tell LangChain if we want to stream generation or not





## Result of **first** question

*Printed text appears  
to be identical to the  
output*

**Good!**

India  
produces  
the  
most  
spices  
,  
contributing  
to  
75  
%  
of  
global  
spice  
production  
.

## Result of **second** question

*The question????*

**???**

*Second group of  
prints appear to be  
the answer*

**Good!**

How  
much  
spice  
does  
India  
produce  
?  
  
India  
produces  
  
1  
,  
525  
,  
000  
metric  
tonnes  
of  
spices  
in  
  
201  
1  
,  
according  
to  
the  
given  
information  
.

## Input

What country produces the most  
spice?

## Conversational Retrieval Chain First message of conversation

*Memory*

## Condense Question Chain

*If memory is empty,  
skip me.*

## Combine Docs Chain

What country produces  
the most spice?

[-.31, .45, .91, ...]

**Retriever**

India is a top  
spice  
producer.

Use the following  
context to  
answer the user's  
question.  
Context: India is a  
top spice  
producer.  
User's question:  
What country  
produces the  
most spice?

**LLM**

**Streaming  
Handler**

India  
produces the  
most spice



## Input

How much?

## Conversational Retrieval Chain After first message of conversation

### Memory

What country produces the most spice?

India produces the most spice

### Condense Question Chain

*Prompt*

Given the following conversation and a follow up question, rephrase the follow up question to be a standalone question, in its original language.

Chat History:  
Human: What country produces the most spice?  
Assistant: India produces the most spice  
Follow Up Input: How much?  
Standalone question:

**LLM**

**Streaming  
Handler**

Exactly how much  
spice does India  
produce?

*Refined Question*

Exactly how much spice  
does India produce?

[-.31, .45, .91, ...]

Retriever

in 2011 india  
produced 1.5m  
tons of spice

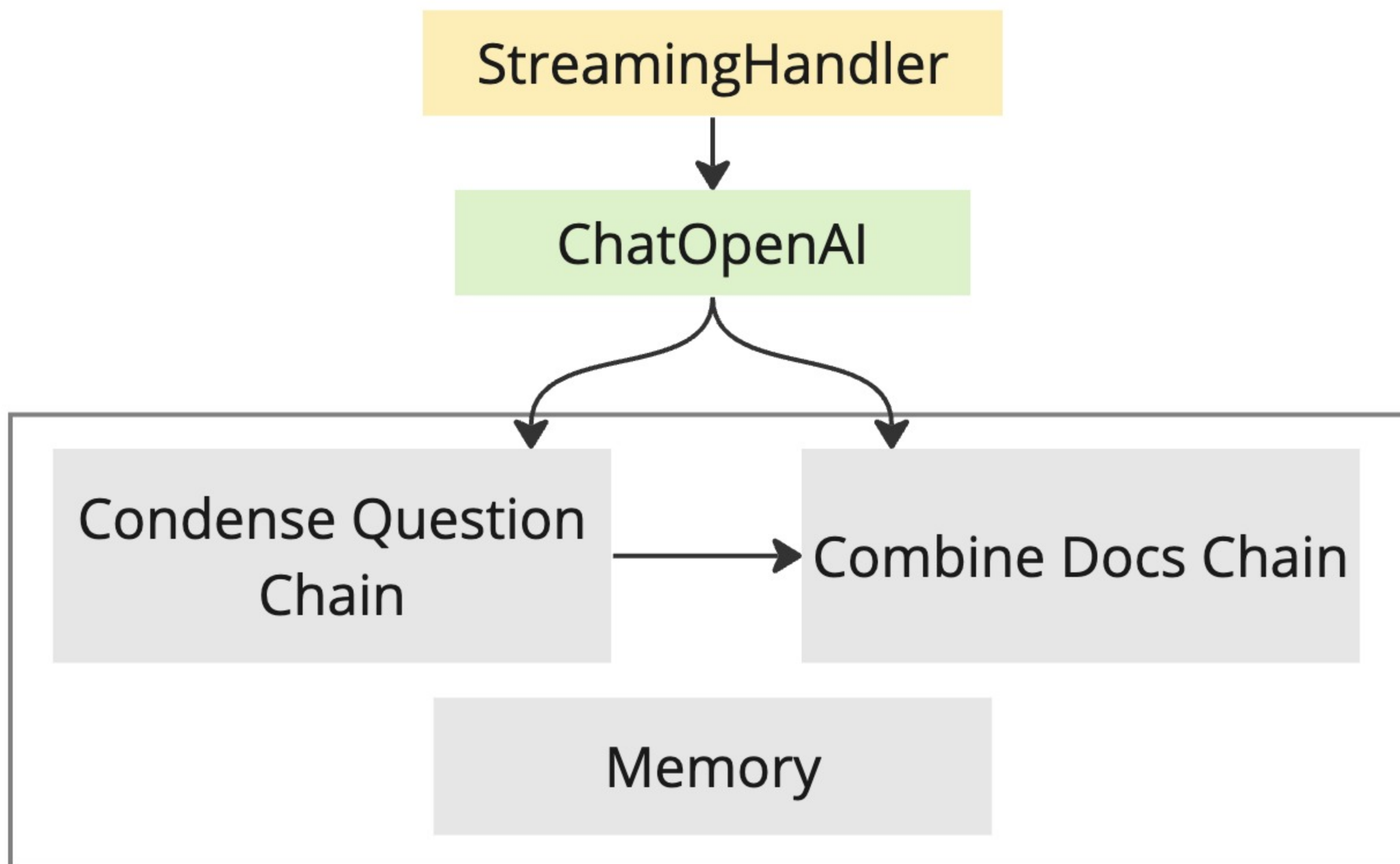
### Combine Docs Chain

Use the  
following  
context to  
answer the  
user's question.  
Context: In 2011,  
India produced  
1.5m tons of  
spice  
User's question:  
Exactly how

**LLM**

**Streaming  
Handler**

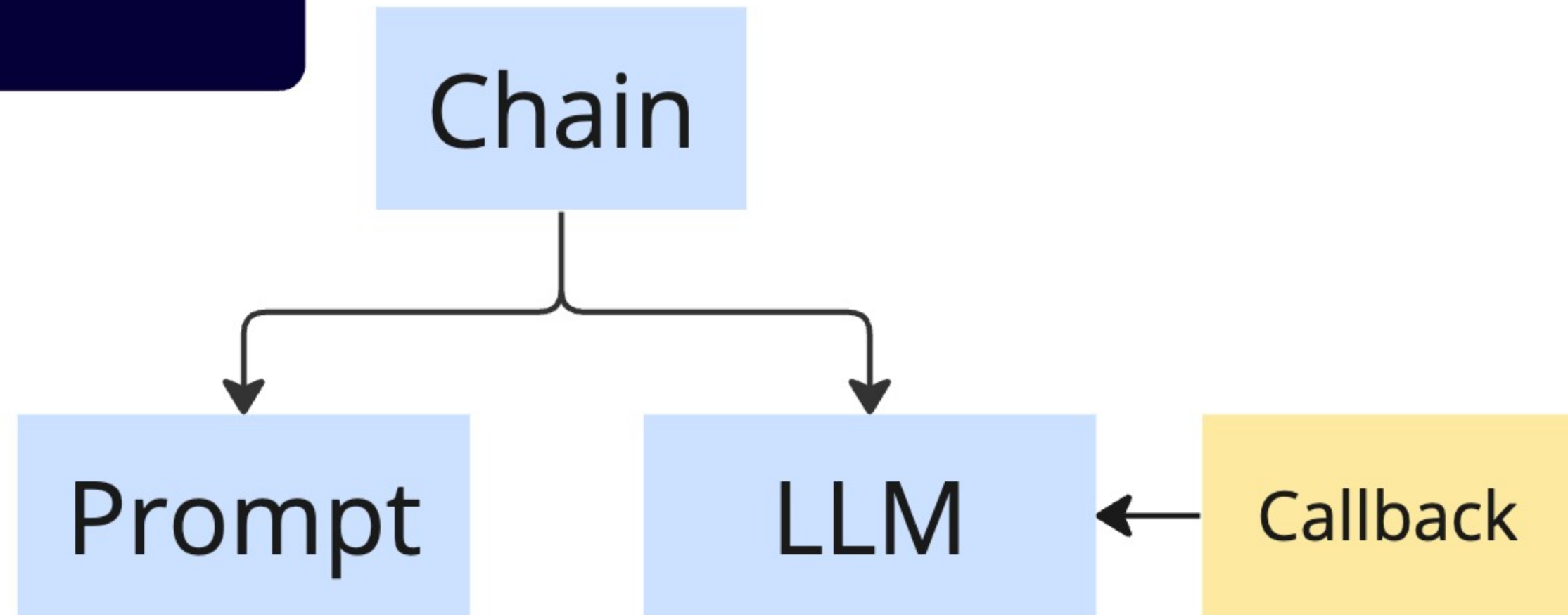
India  
produced  
1.5m tons in  
...





## Callback Applied at Creation

```
1  
2 llm = ChatOpenAI(callbacks=[  
3     MyCallback()  
4 ])  
5 chain = LLMChain(  
6     llm=llm,  
7     prompt=prompt  
8 )
```



## Callback Applied at Call

```
1 result = chain(  
2     "What country...",  
3     callbacks=[MyCallback()]  
4 )
```

