

Building RAGA

NVIDIA®



Building RAGA

NVIDIA®

















Retrieval Augmented Generation



Prerequisites

RAG Agents in Product

• Prior IMM and Chain

Course Objectives

RAG Agents in Product

• Environment+

Building RAGA

NVIDIA®



Last Modified

Name ▾

—

Notebook



/



Filter files by name



Launcher



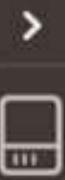
File Settings Help

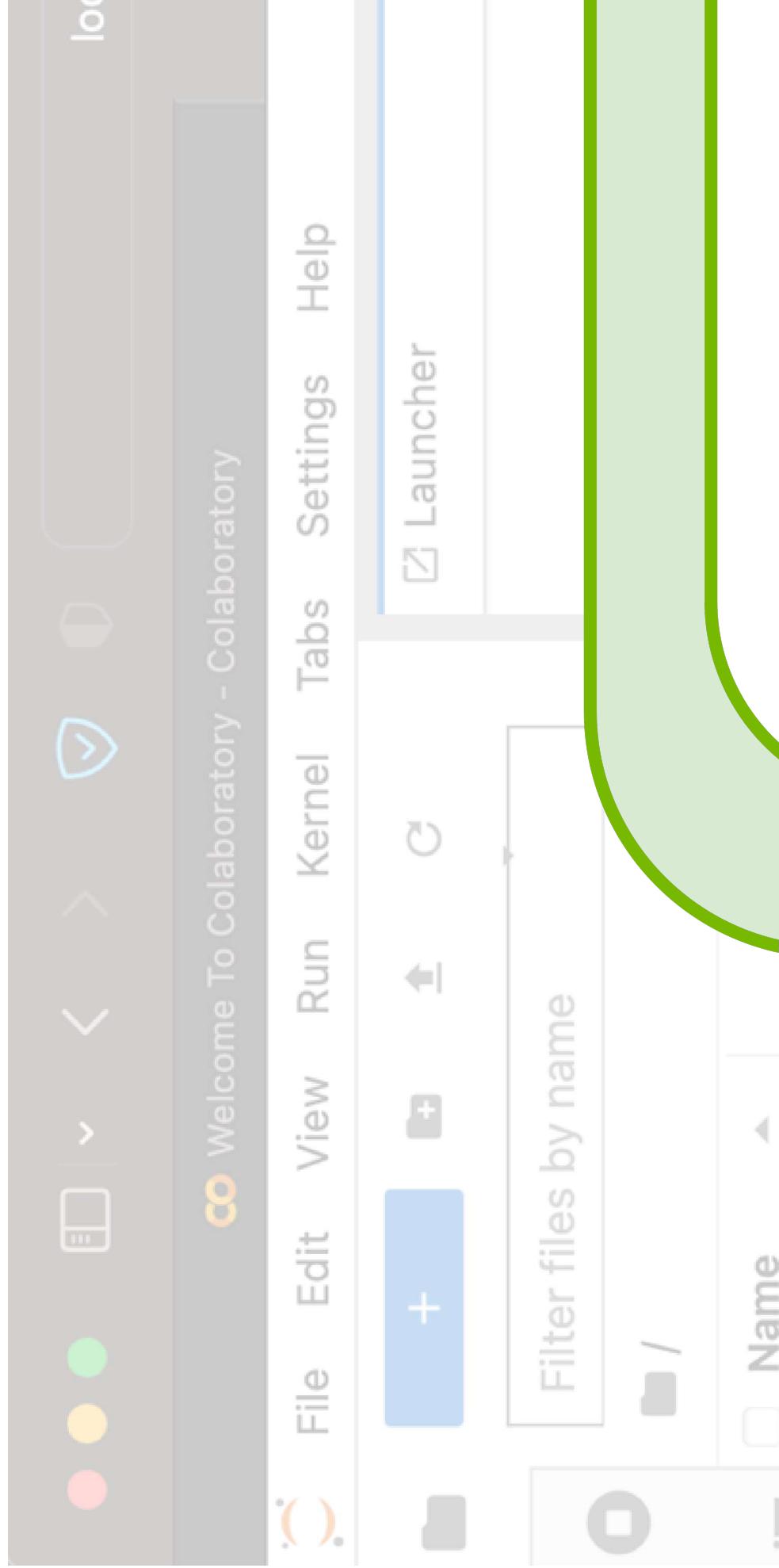


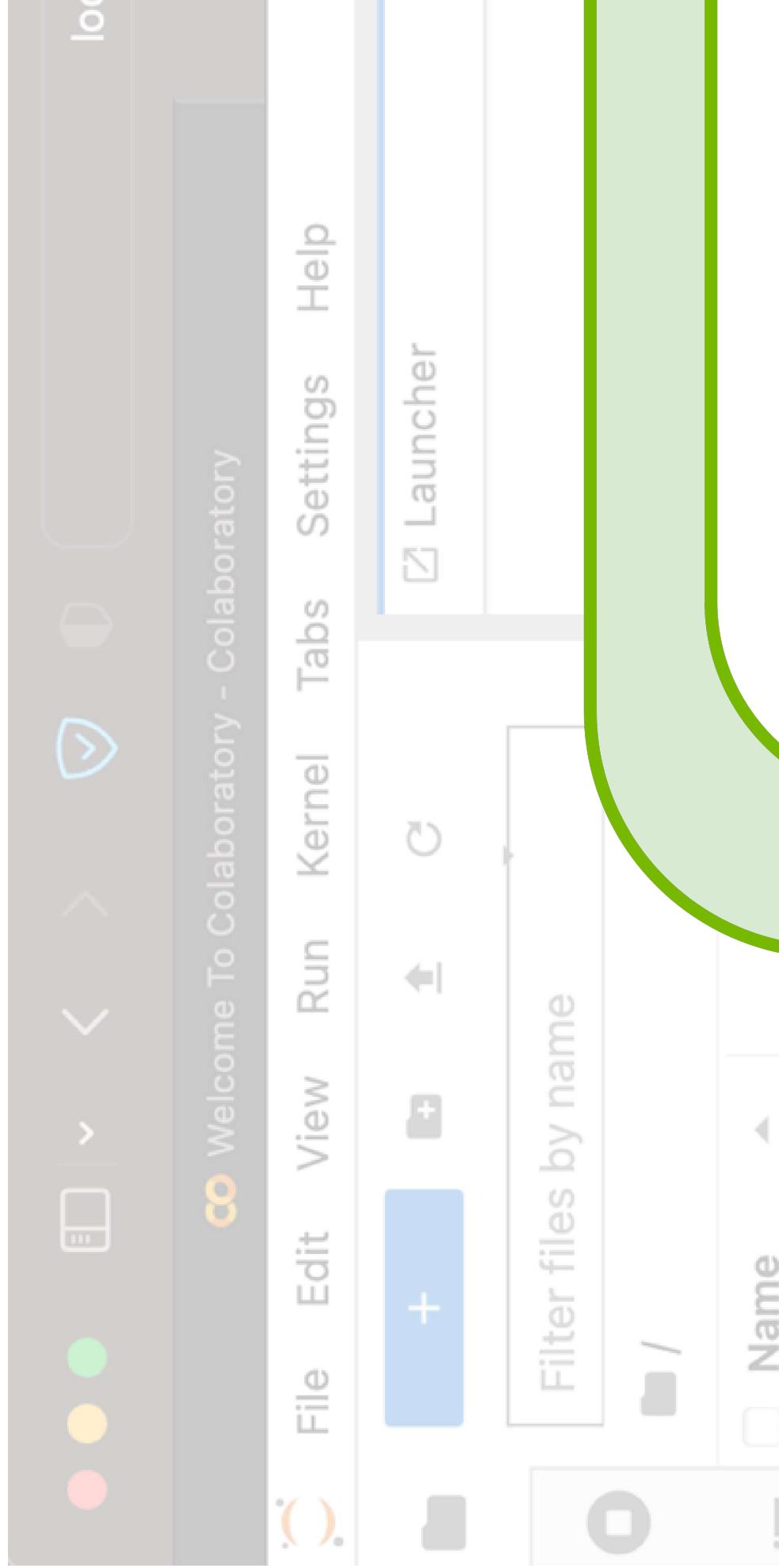
Edit View Run Kernel Tabs

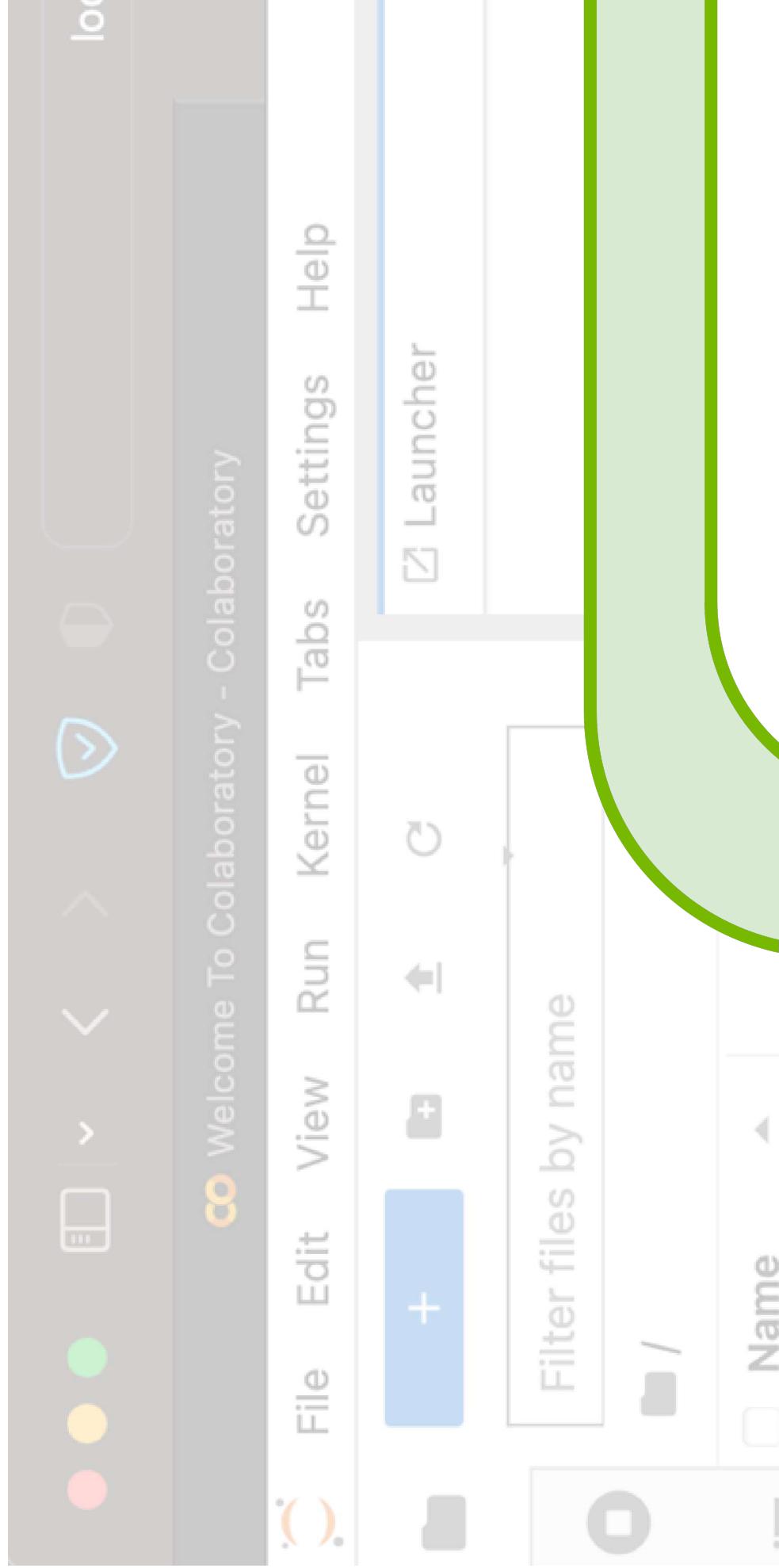
Welcome To Colaboratory - Colaboratory

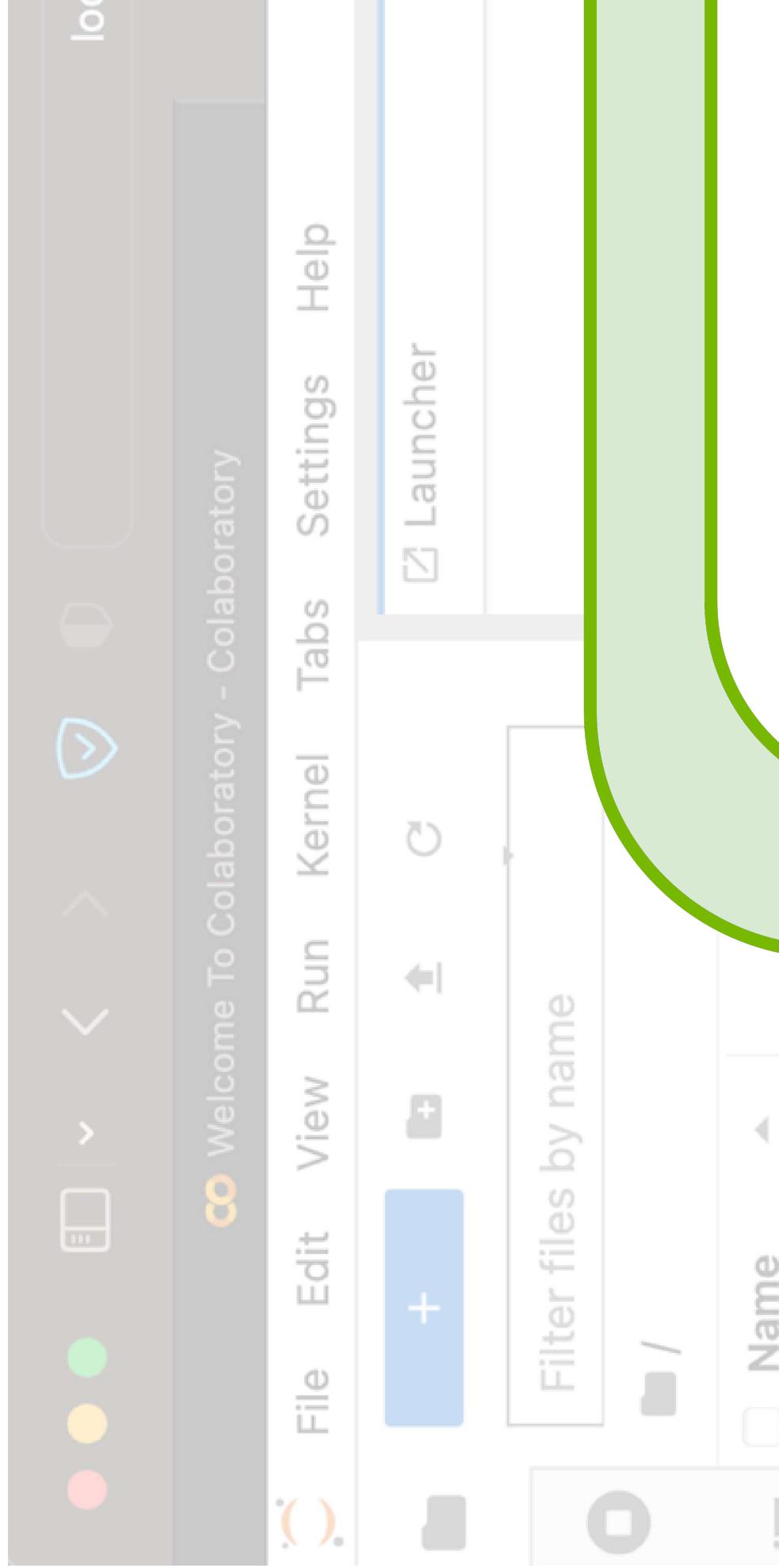
loc

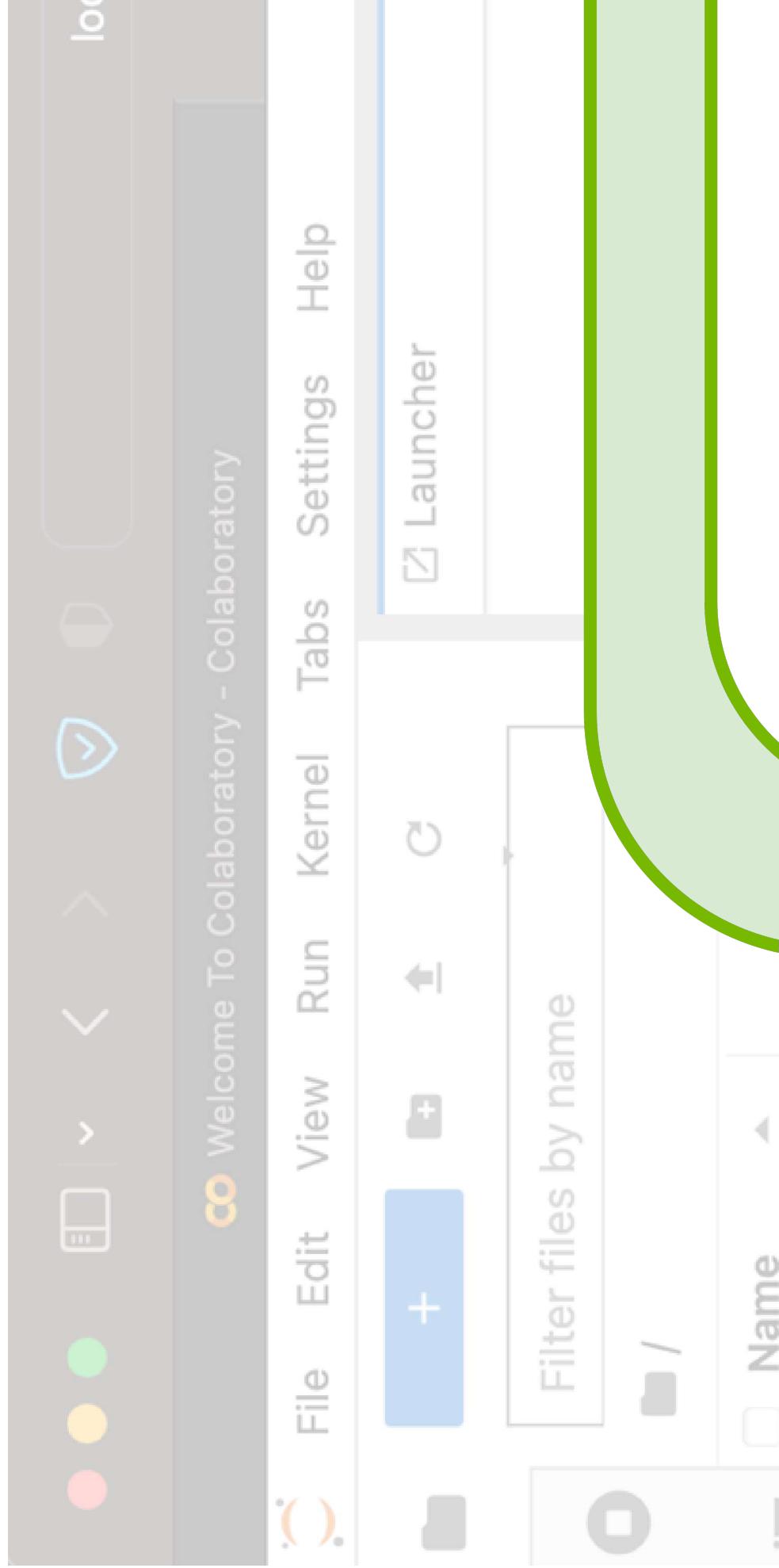


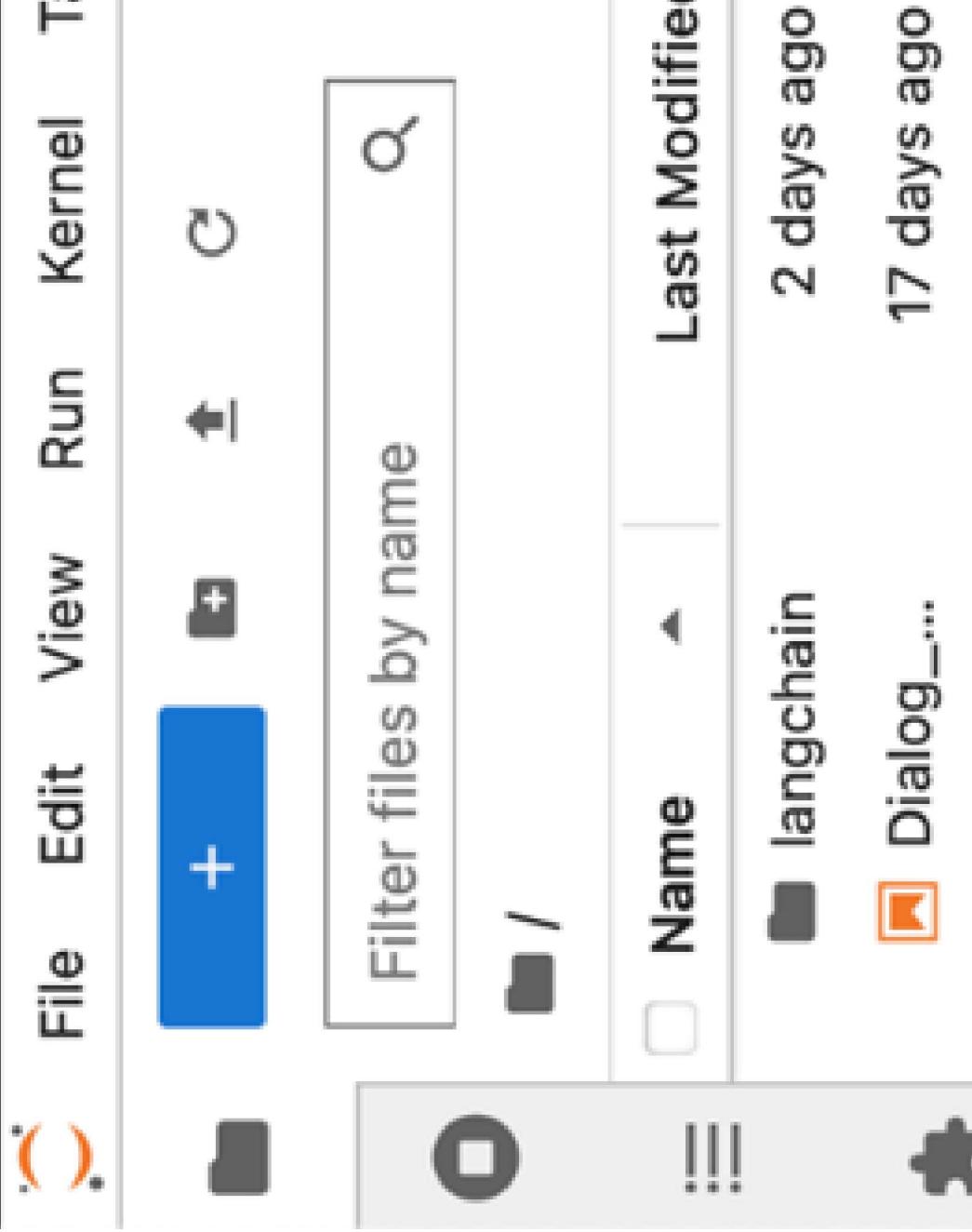


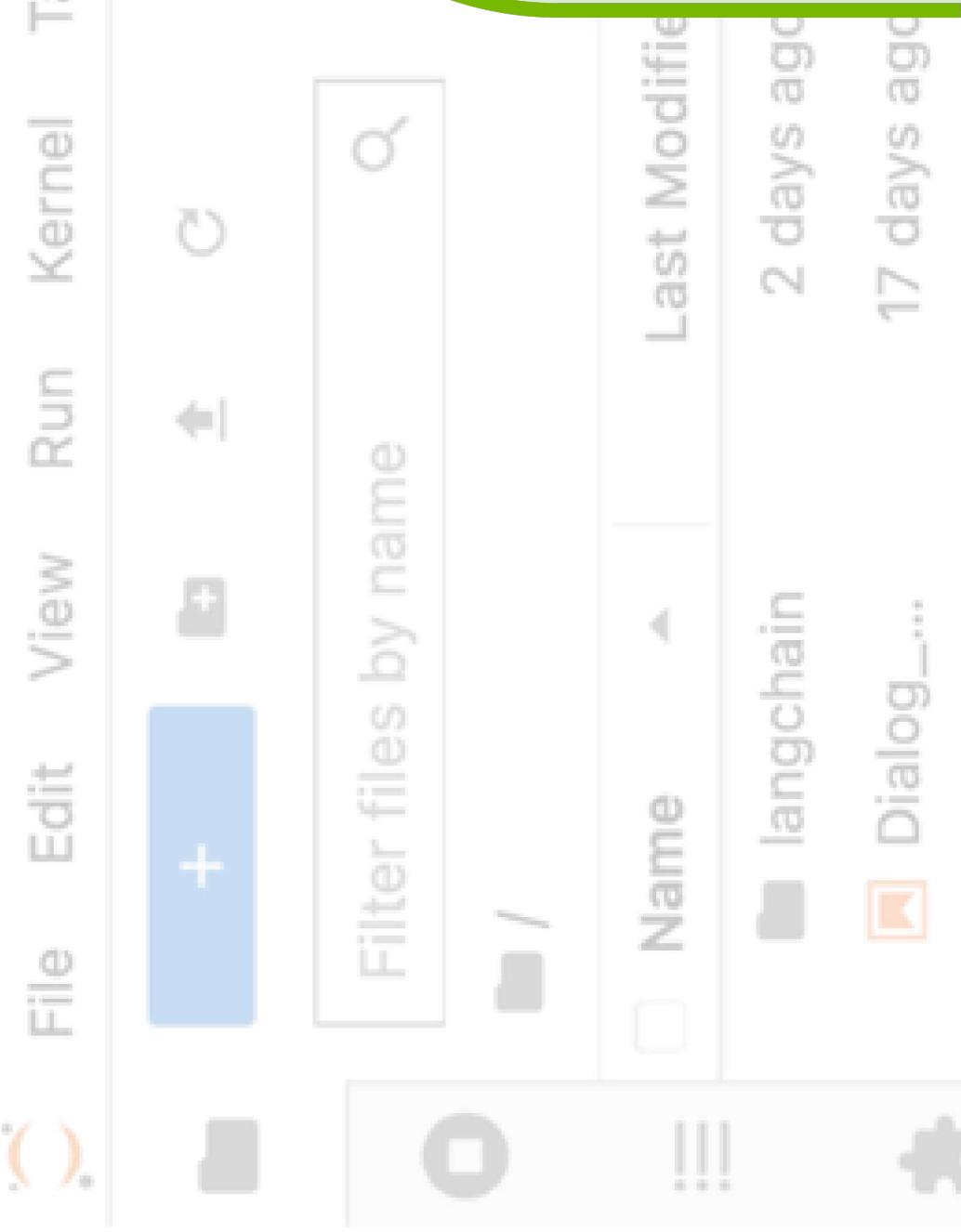


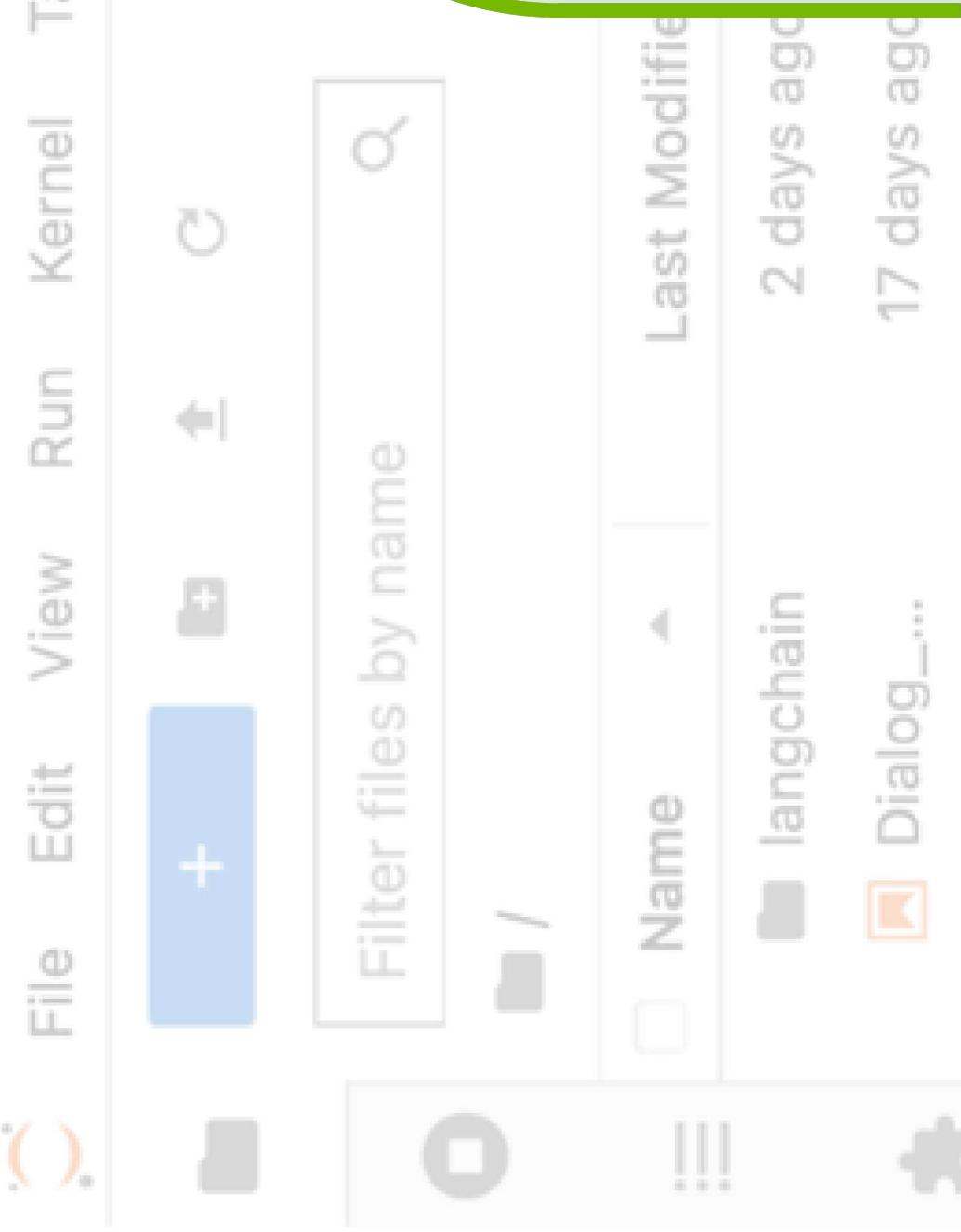
















C

Ho

Schedu

Data

- 1. Allocate Resources**
- 2. Define Services**

Company Database

**Company
Database**

 Chatbot

Hello World!!

You typed: Hello World!!

gradio



Hello World!!

Chatbot

Hello World

Airbnb Map

Diffu

Chatbot Streaming

 Chatbot

 Chatbot

Hello World!!

You typed: Hello World!!

Building RAGA

NVIDIA®



 Chatbot

Hello World!!

You typed: Hello World!!





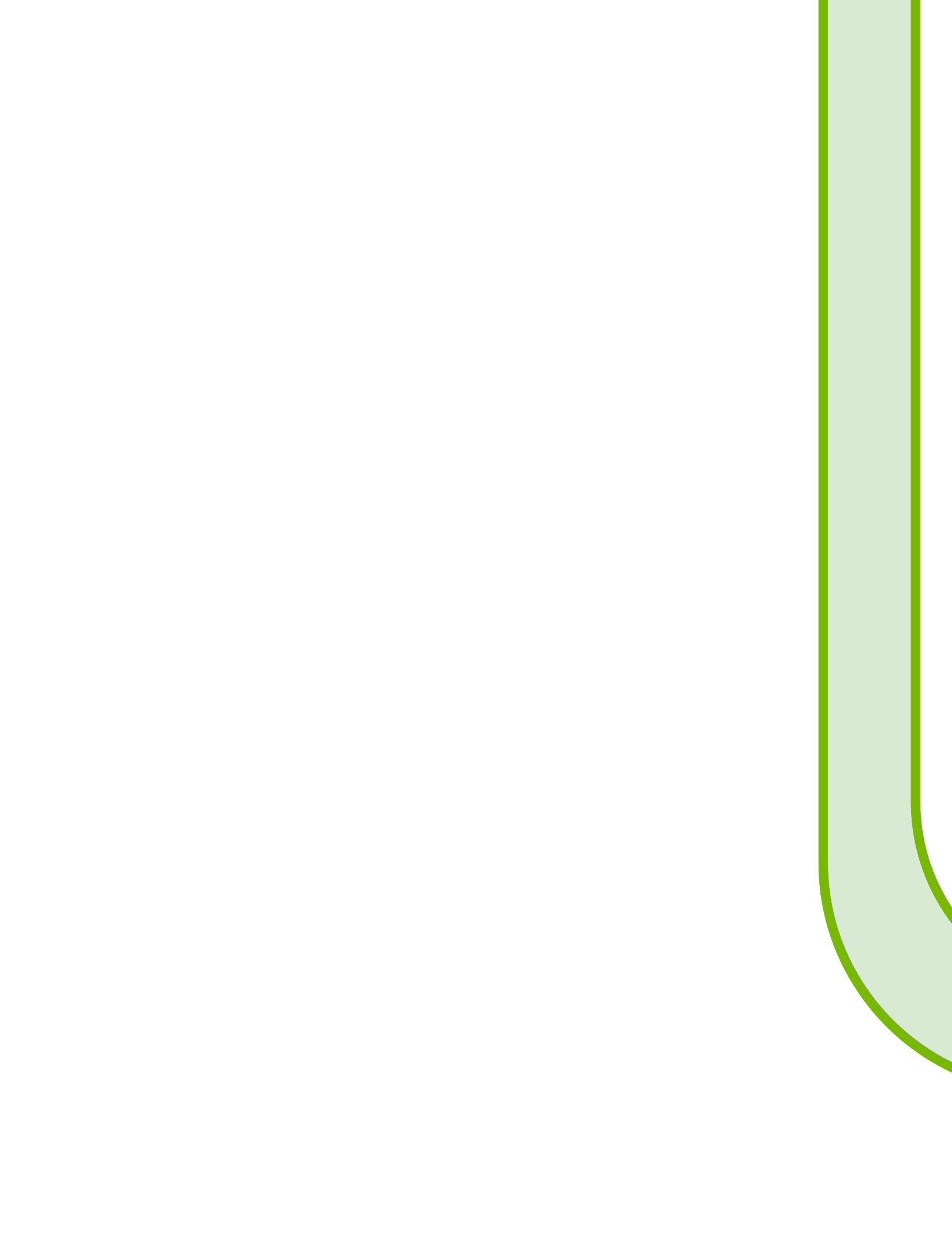


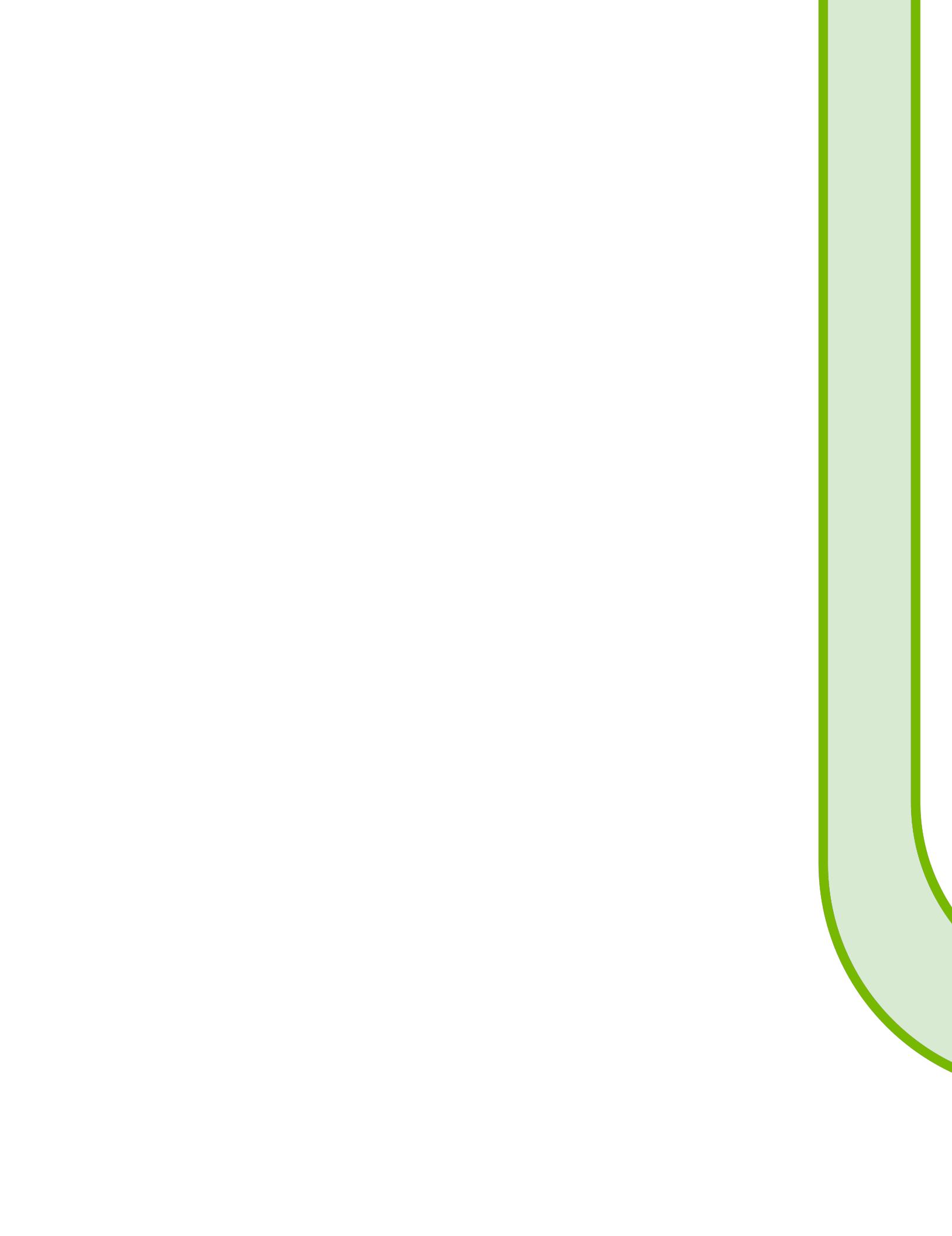


Jupyter Notebook

Jupyter Notebook

Jupyter Notebook








```
  "messages": [ {
```

```
    }
```

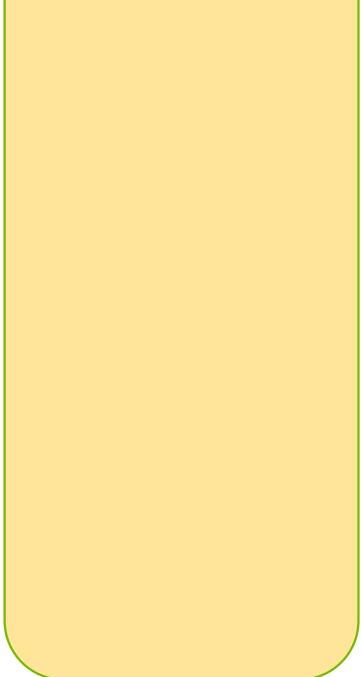
```
  "messages": [ {
```

```
    }
```

Larg

Larg

Larg





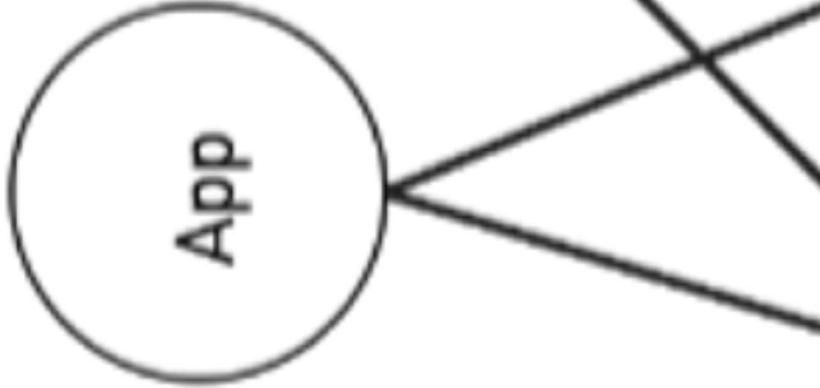


"messages": [{

}

Millions of Apps |

Application Layer



 Chatbot

Hello World!!

You typed: Hello World!!

NV



Search NVII

Top Open Found

The leading open mode

Discover

MODELS

N

NVIDIA.



advanced reasoning

chat

large

An MOEILM that follows instructions

mistral / mixtral-8x

N

From

From



Output
Parsers



Models

Prompts

Example
Selectors

ChatNVIDIA

Example
Selectors

Prompts

Models

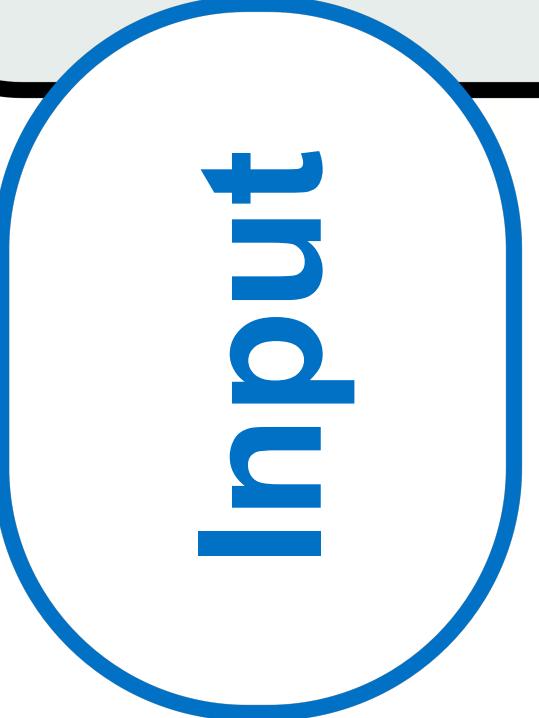
Output
Parsers



Building RAGA

NVIDIA®





Input

Prom

Input

Prom

Input

Prom

Input

OBSI

DEPL

 Chatbot

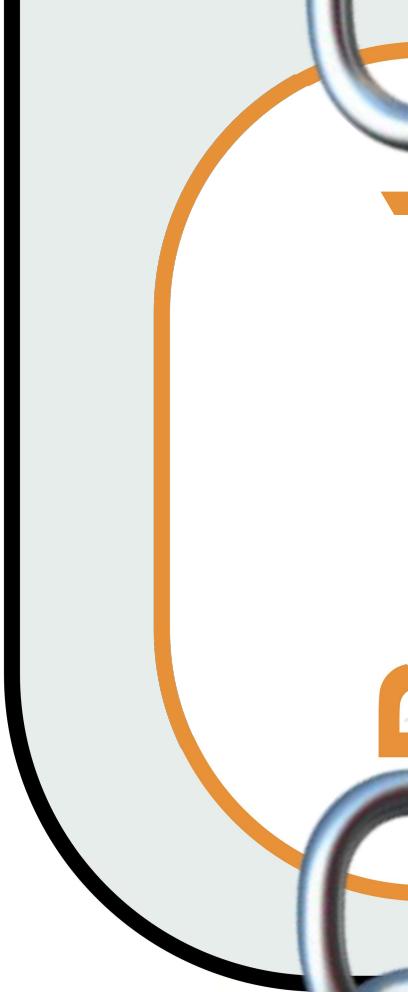
Hello World!!

You typed: Hello World!!

Building RAGA

NVIDIA®





Prompt
Classify
Sentence

Input

```
def fn():  
    ## Initial State
```

$\alpha = 8$

```
F  
def get_initial_state():  
    n = 8  
    fib = [0, 1]  
    return locals()
```

```
R
running_state = { 'n': 8,
                   'n': [7,
                         8,
                         9],
                   'fib': [1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89]}
```

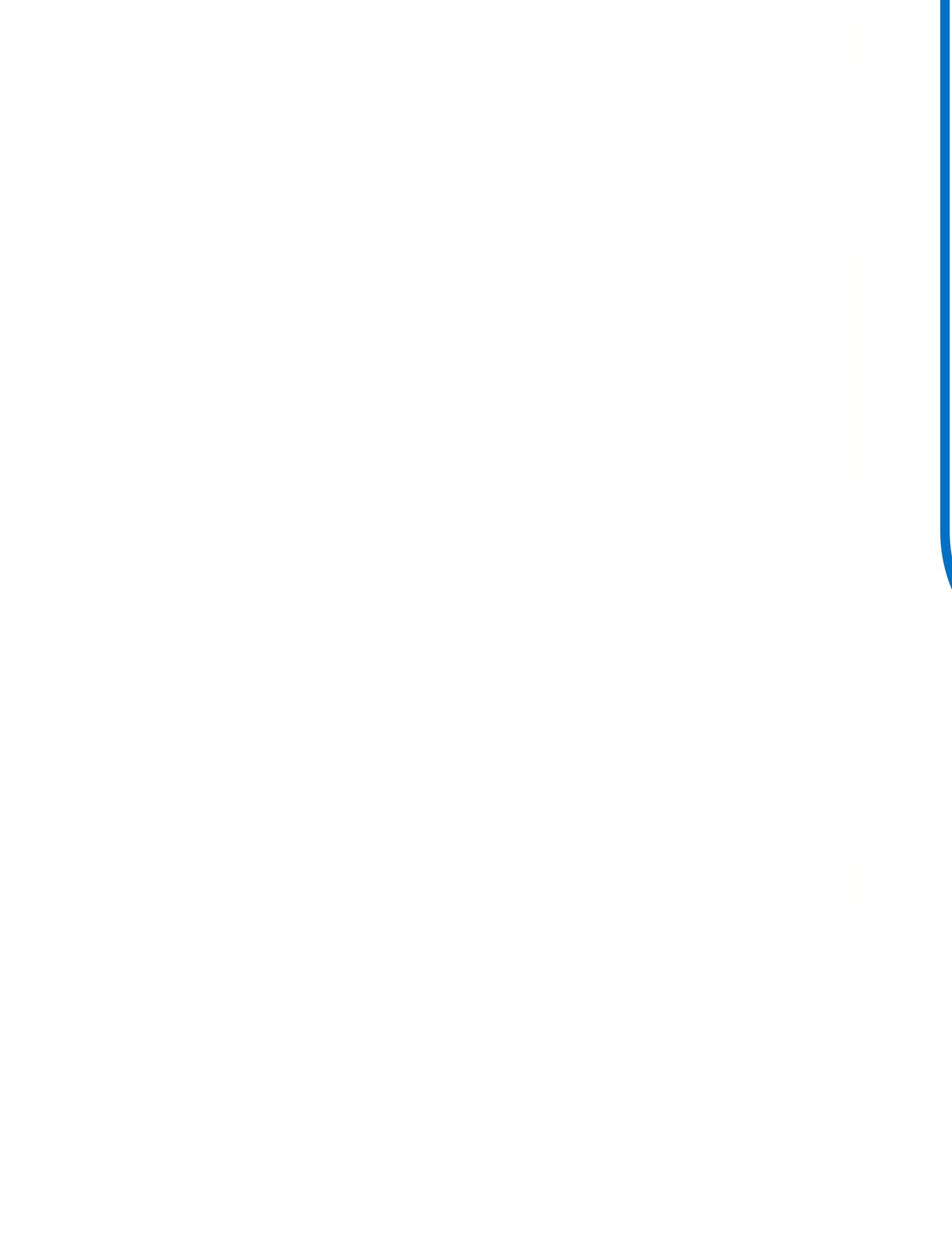
```
F
from langchain_core.runnables
def next_fib_fn(d):
    return d[1] + [d[1], fib
```

```
F
from langchain_core.runnables
def next_fib_fn(d):
    return d[1] + [d[1], fib
```

```
F
from langchain_core.runnables
def next_fib_fn(d):
    return d[1] + [d[1], fib
```

```
from    import re
def    newfile
      
```

```
def fn():  
    ## Initial State  
    n = 8  
    fib = [0, 1]
```



D-a

+|know|

Pro



Pro



Pro



Pro



{

“**first_name**”
“**Last_name**”
“**nickname**”

D-2

+|Know|

D-a

+|know|

Da

+ | know |

Building RAGA

NVIDIA®



Pro



RAG 101: Retrieval-Augmented

Deploying R

Generative AI / LLMs

Generative AI / LLMs

Learn the latest advancements in technology and get hands-on training at GTC 2024. March 18-21.

Learn the latest adva

Technical Blog Search blog Filter

Technical Blog See

NVIDIA DEVELOPER Search blog Filter

NVIDIA DEVELOPER See

Home Blog Forums Docs Downloads Training

NVIDIA DEVELOPER Search blog Filter

NVIDIA DEVELOPER See

Company Database



Company Database

DATA
DRIVEN
DECISIONS

Company Database

Company Database

LLAMA 2: Open Fou

Hugo Peter Albert Amjad Alm Prajwal Bhargava Shruti Bhosale Guillem Cucurull David H. Cynthia Gao Vedanuj Goswami Hakan Inan Marcin Kardas V. Punit Singh Koura Marie

LLAMA 2: Open Foundation and Fine-Tune

Hugo Touvron* Louis Martin† Kevin Stone†
Peter Albert Amjad Almahairi Yasmine Babaei Nikolay Bashlyko
Prajjwal Bhargava Shruti Bhosale Dan Bikel Lukas Blecher Cristian Cai
Guillem Cucurull David Esiobu Jude Fernandes Jeremy Fu Wenyi
Cynthia Gao Vedanuj Goswami Naman Goyal Anthony Hartshorn Saq
Hakan Inan Marcin Kardas Viktor Kerkez Madian Khabsa Isabel Klou
Punit Singh Koura Marie-Anne Lachaux Thibaut Lavril Jenya Lee
Yinghai Lu Yuning Mao Xavier Martinet Todor Mihaylov Pus
Igor Molchanov Yixin Nie Andrew Pouliquen Jeremy Reizenstein Bashi R

LLAMA 2: Open Foundation and Fine-Tuned

Hugo Touvron* Louis Martin† Kevin Stone†
Peter Albert Amjad Almahairi Yasmine Babaei Nikolay Bashlykov S
Prajjwal Bhargava Shruti Bhosale Dan Bikel Lukas Blecher Cristian Canton
Guillem Cucurull David Esiobu Jude Fernandes Jeremy Fu Wenyin F
Cynthia Gao Vedanuj Goswami Naman Goyal Anthony Hartshorn Sagha
Hakan Inan Marcin Kardas Viktor Kerkez Madian Khabsa Isabel Klouma
Punit Singh Koura Marie-Anne Lachaux Thibaut Lavril Jenya Lee Di
Yinghai Lu Yuning Mao Xavier Martinet Todor Mihaylov Pushk
Igor Molybog Yixin Nie Andrew Poultou Jeremy Reizenstein Rashi Run
Alan Schelten Ruan Silva Eric Michael Smith Ranjan Subramanian Xiaocing

How does Flying t
work according

Runnable

```
from langchain_openai import ChatOpenAI
from langchain_core.messages import Message
from langgraph.graph import END, Message
```

Runnable

Building RAGA

NVIDIA®



How does Flying t
work according

Pro



Pro



(Common)

Autoregression

Deep

(Common)

Autoregression

Deep



High-performance com

High-performance com

High-performance com

Bi-En

Symmetry



f(query,

High-performance com

Building RAGA

NVIDIA®



High-performance com



**Irrelevant
Questions**

Irrelevant

**Illegal
Topics**

Irrelevant

**Illegal
Topics**

Irrelevant

**Illegal
Topics**



**Irrelevant
Questions**

Building RAGA

NVIDIA®



High-performance com

Who is the



Vector Stores

```
embedder = NVIDIA_Embbedder
```

```
from Langchain_nvidia_ai  
from Langchain_vectorstores
```

```
embedder = NVIDIA_Embbedder
```

```
from Langchain_nvidia_ai  
from Langchain_vectorstores
```

```
embedder = NVIDIA_Embbedder
```

```
from Langchain_nvidia import  
from Langchain.vectorstores
```

```
embedder = NVIDIA_Embbedder
```

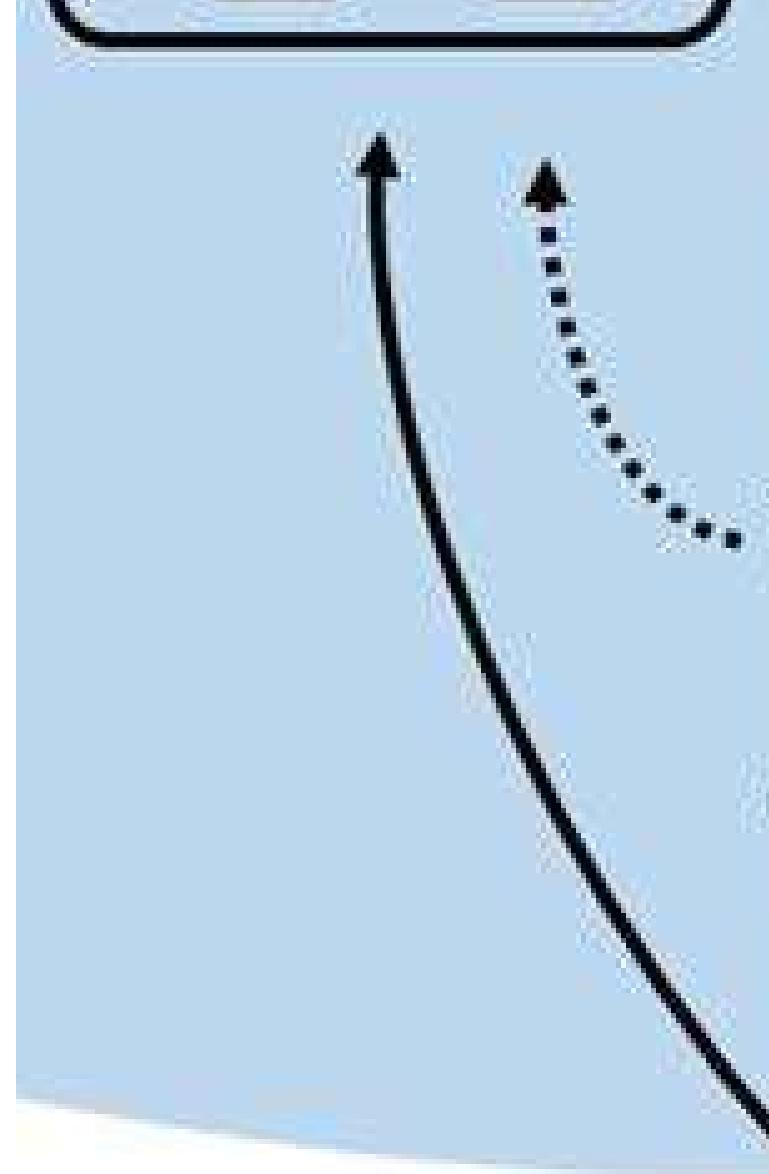
```
from Langchain_nvidia_ai  
from Langchain_vectorstores
```


Pron



Pron





Loca

Local Host

Milvus

Batch size 10, H100 (SXM)

DEEP-100M vector search throughput

Hello Jane!
How are you?

Hello! My
name is Jane



How does RAG work?



Building RAGA

NVIDIA®



RAG Pipeline

How do

RAG Pipeline

How do

vDB

vDB

vDB

vDB

Gel

VDB

+ Embed

U.S.

F Q

NLI

54 ✓ NLI

55

Na

56

57

58

59

Con

Joh



LLM

Prompt
Ask some
questions

vDB

Remote



Congratulation

NVIDIA®

