# Jupyter-AI & Open LLMs in large science classes

ESPM
157

**Carl Boettiger**
**UC Berkeley**

- **122** students
- ~ **2/3** bio sciences
- **1/3** data science
- Active learning classroom



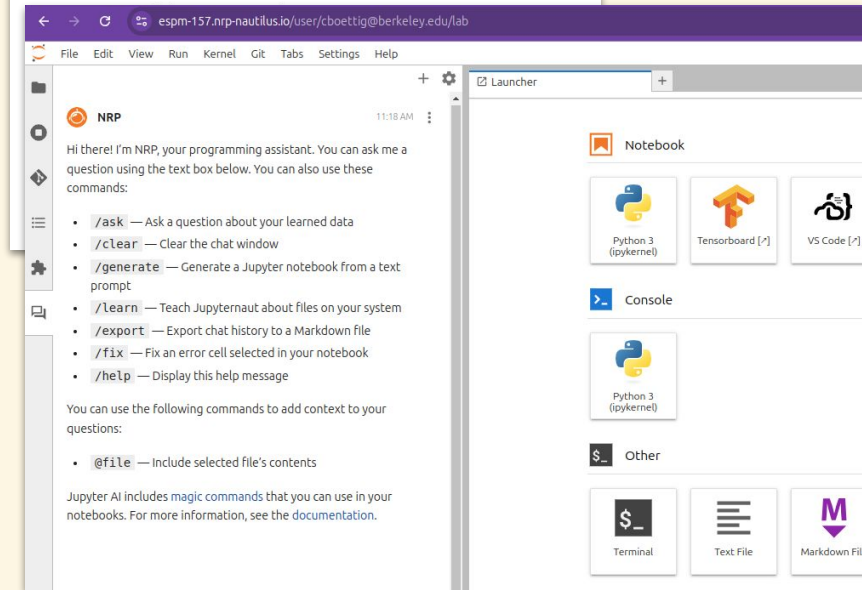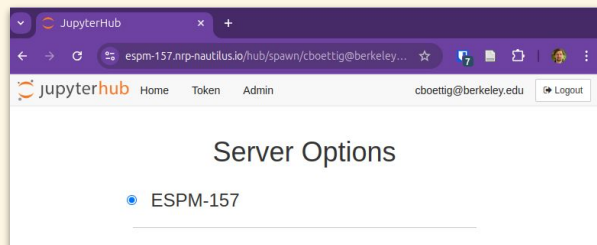# ESPM-157: Data Science for Global Change Ecology

# Motivations

*LLMs are here. Am I preparing my students for this future?*

- Can we go further, faster by **coding** with LLM assistants?
- What are **environmental** implications of LLMs?
- What are the **ethical** implications?
- Can we make LLM use **safer** & more **reliable**?

# Computing Resources

**NRP NATIONAL RESEARCH PLATFORM**

**nvidia.com/gpu**: "1"  x 122 students!



**Environment.yml Docker image**

**k8s Helm chart**

# Challenges with respect to AI

- Student data privacy?
- Energy footprint?
- Ethical concerns:
  - copyright,
  - labor practices,
  - biases

*those who concern themselves with such things as "**the computer and its social impacts**" but who fail to **look behind technical things** to notice the **social circumstances** of their development, deployment, and use.*

*This view provides an **antidote** to naive **technological determinism** — the idea that technology develops as the sole result of an internal dynamic, and then, unmediated by any other influence, molds society to fit its patterns. Those who have not recognized the ways in which technologies are shaped by social and economic forces have not gotten very far.*

- *Langdon Winner, 1980*

ARTICLE
https://doi.org/10.1057/s41599-024-02720-3  OPEN

AI chatbots contribute to global conservation injustices

Danilo Urzedo[1&], Zarrin Tasnim Sworna[1], Andrew J. Hoskins[2] & Cathy J. Robinson[1,3]

collection and translation of environmental evidence that could be used to inform planetary conservation plans and strategies. Yet, the consequences of chatbot-generated conservation

AI and the Copyright Liability Overhang: A Brief Summary of the Current State of AI-Related Copyright Litigation
Article  April 2, 2024  6 minutes

# Local, open models:

```
$ ollama serve
$ ollama pull nomic-embed-text
```



File   Edit   View   Run   Kernel   Git   Tabs   Settings   Help

Language model

Completion model
NRP :: llama3

Click here for more details on NRP

Embedding model

Embedding model
NRP :: embed-mistral

| ERNIE-Bot :: ERNIE-Bot-4 |
| GPT4All Embeddings :: all-MiniLM-L6-v2-f16 |
| MistralAI :: mistral-embed |
| NRP :: embed-mistral |
| Ollama :: nomic-embed-text |
| Ollama :: mxbai-embed-large |
| Ollama :: all-minilm |

# **CO2** footprint data, not hype:



## Global

| Energy consumed | Emissions produced |
| --- | --- |
| **222.06** | **94.28** |
| kWh | Kg. Eq. CO2 |

CODE CARBON
Track and reduce CO2 emissions from your computing

01/01/2020 → 01/02/2022



OLMo-7B
OLMo-7B-Instruct
AI2

# Researcher developed + hosted



```
! llm.yml
  6    spec:
 15      template:
 20        spec:
 21          affinity:
 30          containers:
 31          - args:
 32            - -m
 33            - vllm.entrypoints.openai.api_server
 34            - --port
 35            - "5000"
 36            - --host
 37            - 0.0.0.0
 38            - --download-dir
 39            - /workspace/.cache/huggingface/hub
 40            - --model
 41            - gorilla-llm/gorilla-openfunctions-v2
 42            - --tensor-parallel-size
 43            - "2"
 44            - --trust-remote-code
 45            - --enable-auto-tool-choice
 46            - --tool-call-parser
 47            - llama3_json
 48            image: vllm/vllm-openai:v0.6.3
 49            imagePullPolicy: IfNotPresent
 50            name: gorilla-openfunctions-v2
 51            command:
 52            - python3
 53            resources:
 54              limits:
 55                cpu: "5"
 56                memory: 36Gi
 57                nvidia.com/rtxa6000: "2"
 58              requests:
 59                cpu: "1"
 60                memory: 36Gi
 61                nvidia.com/rtxa6000: "2"
 62
```

> base_url = "https://llm.nrp-nautilus.io/"

# Challenges with respect to AI

- LLMs **don't**

/learn docs
/ask ...

# RAG & TAG – Building LLM Agents

```python
st.title("SQL demo")

## dockerized streamlit app wants to read from c…
api_key = os.getenv("LITELLM_KEY")
if api_key is None:
    api_key = st.secrets["LITELLM_KEY"]


parquet = st.text_input("parquet file", "https:/…


eng = sqlalchemy.create_engine("duckdb:///:memory:…
con = ibis.duckdb.from_connection(eng.raw_connecti…
tbl = con.read_parquet(parquet, "mydata")


# langchain can also talk to this connection and s…
db = SQLDatabase(eng, view_support=True)


# Build the template for system prompt
template = '''
You are a {dialect} expert. Given an input questio…
Always return all columns from a query (select *)…
Wrap each column name in double quotes (") to deno…
Pay attention to use only the column names you can…
Be careful to not query for columns that do not ex…
Also, pay attention to which column is in which ta…
Pay attention to use today() function to get the c…
Respond with only the SQL query to run.  Do not re…
Only use the following tables:
{table_info}
Question: {input}
'''

from langchain.core.prompts import PromptTemplate
```

app

chat

rag

sql

## SQL demo

parquet file:

https://espm-157-f24.github.io/spatial-carl-amanda-tyler/new_haven_stats.parquet

SELECT grade, AVG(ndvi) AS mean_ndvi FROM mydata GROUP BY grade;

|   | grade | mean_ndvi |
|---|-------|-----------|
| 0 | A     | 0.7308    |
| 1 | C     | 0.5472    |
| 2 | D     | 0.4866    |
| 3 | B     | 0.6093    |

What is the mean ndvi by grade?

# Project Contributions & Future engagement

- Teach **safe LLM** use, not **AI abstinence**
- **Open models** mitigate some risks (energy, privacy)
- NRP as platform for collaborative innovation!
  - **Shared** hardware, software, models