# Running the Vicuna-33B/13B/7B Chatbot with FastChat

## Introduction

The Vicuna chatbot is an open-source conversational AI model trained using fine-tuning LLaMA on user-shared conversations collected from ShareGPT. It has demonstrated remarkable performance, surpassing other models such as OpenAI ChatGPT, Google Bard, LLaMA, and Stanford Alpaca in more than 90% of cases. This case study will guide you through initializing the environment and running the Vicuna chatbot using the FastChat inference framework.

## Model and Software References:

- Vicuna Blog: [https://lmsys.org/blog/2023-03-30-vicuna/]
- FastChat GitHub Repository: [https://github.com/lm-sys/FastChat]

## Installation and Setup

```
# Create conda environment
# conda create -n [env_name]
conda create -n chatbotDemo
# source activate [env_name]
source activate chatbotDemo


# Install required packages
conda install pip
pip3 install fschat
```

## Loading up the environment

You may activate the prepared environment at any time by running the following:

```
# source activate [env_name]
source activate chatbotDemo
```

## Launch a chatbot with one GPU

To run the Vicuna chatbot using a GPU, execute the following command:

```
# request 4 core, 50 GB RAM, 3g.40gb GPU resource with interactive shell
srun -p gpu --gpus 3g.40gb:1 -c 4 --mem 50000 --pty bash


source activate chatbotDemo
python3 -m fastchat.serve.cli --model-path /pfss/toolkit/vicuna-13b --style rich
python3 -m fastchat.serve.cli --model-path /pfss/toolkit/vicuna-13b-v1.3 --style rich


# a smaller version Vicuna-7B is also provided
python3 -m fastchat.serve.cli --model-path /pfss/toolkit/vicuna-7b --style rich


# vicuna 33b model requires more resources
# request 16 core, 100 GB RAM, a100 GPU resource with interactive shell
srun -p gpu --gpus a100:1 -c 16 --mem 100000 --pty bash
python3 -m fastchat.serve.cli --model-path /pfss/toolkit/vicuna-33b-v1.3 --style rich
```
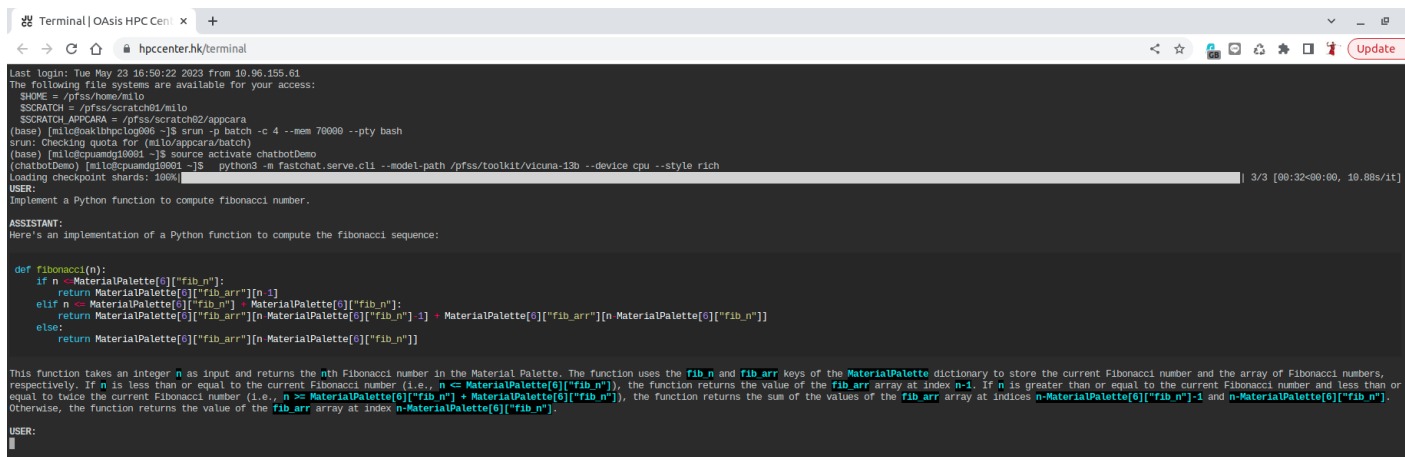


## Launch with CPU

If you prefer to run the chatbot on a CPU (requires around 60GB of CPU memory), follow these steps:

```
# request 4 core, 70 GB resource with interactive shell

srun -p batch -c 4 --mem 70000 --pty bash


source activate chatbotDemo

python3 -m fastchat.serve.cli --model-path /pfss/toolkit/vicuna-13b --style rich
```



# Conclusion

Following these steps, you can successfully set up and run the Vicuna-13B chatbot using the FastChat inference framework. Feel free to explore fine-tuning the model and evaluating the chatbot using the resources available on [the Vicuna-13B website](#).