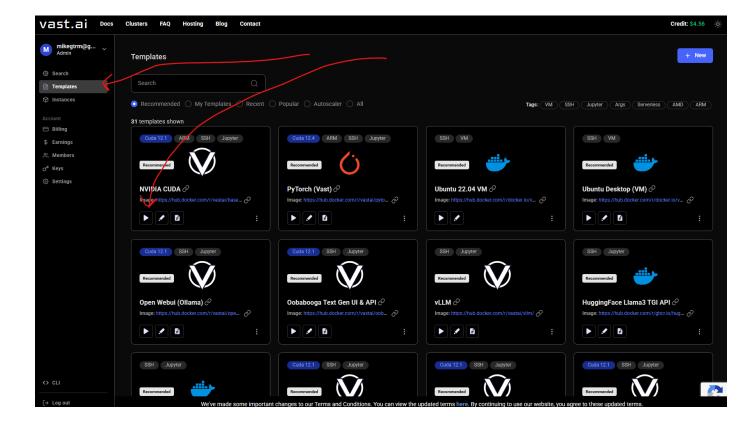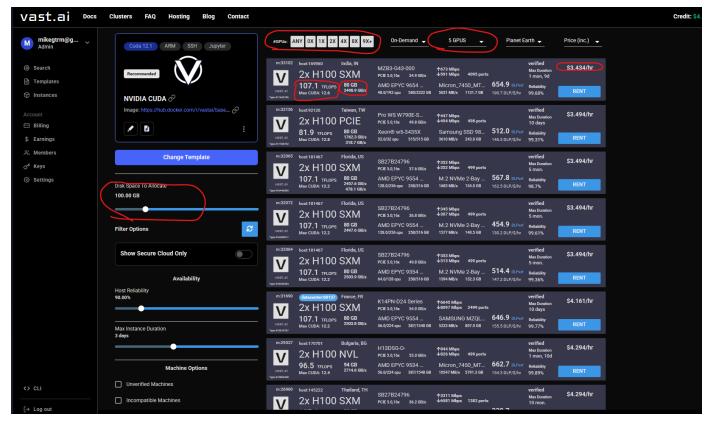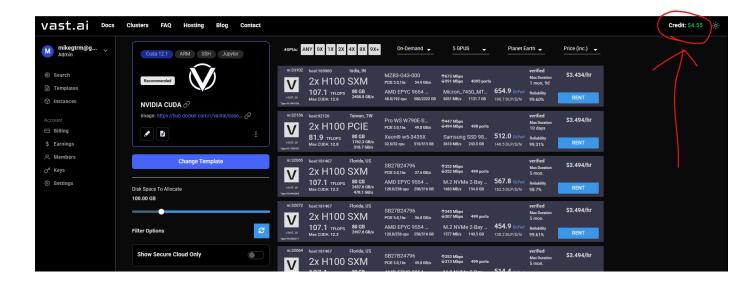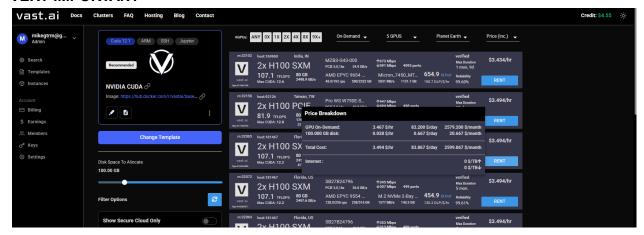# READ ALL BEFORE RENTING
Renting GPUs

- You've got a couple good settings here, make sure you have at least 100gb to allocate (this is disk space) so the space our dataset takes.
- The white buttons let you select how many GPUs.
- The GPU lets you select which GPUs you want (RTX, H-Series, A-Series, etc)
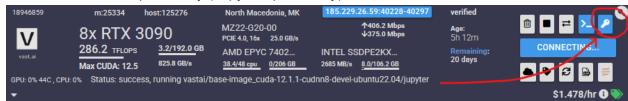- The TFLOPs is floating point operations, higher is typically better
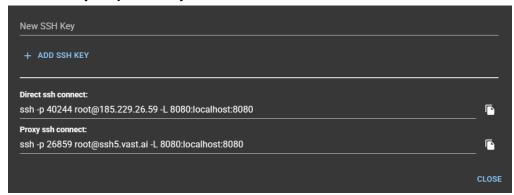
**VERY IMPORTANT**



When you hover over the rent button you get to see the cost breakdown. **Sometimes** you can find some cheap GPUs per hour, but the internet price is like $40 for every Terabyte you download. This typically is only high for GPU providers outside of the US

**AFTER YOU CLICK RENT**

- Create a ssh public key (copy the public key)



- Add your **public** key



- Copy over the ssh command and ssh into it

**AFTER YOU ARE IN THE SYSTEM**

- git clone https://github.com/Michaelgathara/GPT
- cd GPT

- curl -LsSf https://astral.sh/uv/install.sh | sh
- source $HOME/.local/bin/env bash
- uv sync
- source .venv/bin/activate
- uv pip install flash-attn --no-build-isolation
- cd models
- nohup python3 -u gpt_train_script.py > train.log 2>&1 &
- bash print_res.sh