**Setup: AWS EC2 GPU Instance (g6.xlarge)**

For cost-efficiency, I used a g6.xlarge single-GPU instance.

**Launch Steps:**

1. **EC2 Console** - Launch new instance
2. **AMI Selection** - Use "Deep Learning Base AMI with Single CUDA (Amazon Linux 2023) 20260120"
   - Alternatively, choose any base OS and install drivers manually
   - Note: Deep learning Ubuntu 22.X variant had NVIDIA driver detection issues. So this did not work for me
3. **Instance Type** - Select g6.xlarge
   - **Important**: Request vCPU quota increase beforehand (approval may take time; good luck if you account is new)
4. **SSH Access** - Configure key pair (.pem file)
5. **Network** - Default settings; customized security group to allow SSH from my public IP
6. **Storage** - gp3, 125 GB root volume

**Verification:**

After connecting via SSH, run `nvidia-smi` to confirm GPU availability.



In this instance, python was already installed.