

Specialist ML/AI Assignment

AI ML specialist Customer Solution Architect

Thank you for applying for the AI ML specialist Customer Solution Architect role! As your take home assignment, you are asked to complete the following exercises before your demo day.

Please note:

Exercise 1: You are expected to perform a short presentation during the demo day.

Exercise 2: Only code example is required.

Exercise 1: LLM fine-tuning

You have just received some news about a potential customer. Before committing to reserve some GPU capacity in Nebius (they are looking to reserve 512 H100 GPUs for initial duration of 6 months), they would like to do some testing of our platform during a PoC Proof-of-Concept) stage.

About the client: a small VC-funded startup (20 headcount), they work on process automation with AI agents.

You had a brief call with the team on the customer's side who will be running the PoC. They are mostly ML Engineers without extensive Cloud/Infrastructure expertise, and they are looking to do an array of tests including an end-to-end fine-tuning of an open-source LLM to perform function calling.

Your objective is to support this team during the PoC. For this, you will need to prepare an end-to-end example of multi-node fine tuning of an LLM on the capacity allocated for the PoC:

16 H100 GPU cards

2TB SSD network disk

2TB SSD shared filesystem

. The example should utilize the provided PoC capacity efficiently. The choices of scheduler, framework, storage type are up to you. Code example is required.

During the demo day, you will be asked to present your example to the customer and explain your technical choices. You should also provide necessary documentation (how to reproduce and monitor the fine-tuning example).

Useful links:

Nebius CSA Solution Library:  [GitHub - nebius/nebius-solution-library](#)

K8s solution:  [nebius-solution-library/k8s-training at main · nebius/nebius-solution-library](#)

Slurm solution:  [nebius-solution-library/soperator at main · nebius/nebius-solution-library](#)

Nebius AI Cloud documentation: <https://docs.nebius.com/>

Slurm operator for k8s Soperator):  [Soperator](#)

Exercise 2: Solution for testing GPU clusters

Your task is to prepare an internal solution library example. This example is a GitHub repository with code allowing to build a container which runs distributed training of an ML model which may be used in acceptance testing of GPU clusters by Cloud/Infrastructure engineers. The implementation you prepare should be portable and lightweight, using only open-source models/datasets.

Tips:

Use `nvcr.io/nvidia/pytorch:24.07-py3` as base image;

The repository should contain a CI script for building the container and running a simple test;

If tests are successful, the image should be pushed to `ghcr.io`.