

NLP Tasks



Dialogue Systems



Recommender Systems



APPLICATIONS

Compare cats and dogs.

What is the capital of China?

...

Find rhyming word for the word a river.



data collection



predication & profiling



SLO

...

Output length

resource profiler

...

Q3

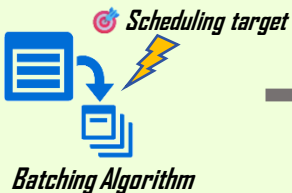
Q2

Q1

Prompt: Compare cats and dogs.

SLO: 200 ms

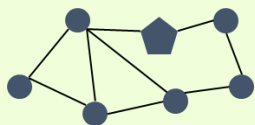
Output length: 20 tokens



Batching Algorithm



batch scheduler



Network topology



LLMs



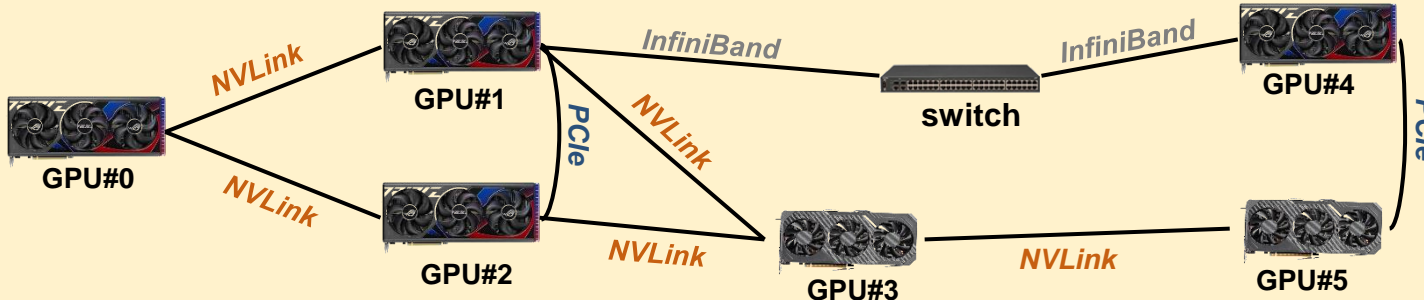
Deployment algorithm



layer	device
layer 1	cuda 0
layer 2	cuda 0
...	...
layer n	cuda m

Config: device map

LLM deployer



GPU-BASED CLUSTER

UELLM