

# SpotServe: Serving Generative Large Language Models on Preemptible Instances

Xupeng Miao, Chunan Shi, Jiangfei Duan, Xiaoli Xi,  
Dahua Lin, Bin Cui, and Zhihao Jia.

**Presenter : Yi Liu**



# Author

## Research Direction:

- machine learning systems
- data management
- distributed computing

## Selected Publications:

- SpotServe(ASPLOS'2024)
- SpecInfer(ASPLOS'2024)
- Galvatron(VLDB 2023)



Xupeng Miao

<https://hsword.github.io/>



# Outline

- Background
- Design
- Evaluation
- Thinking

# Background:The Scale Of LLM

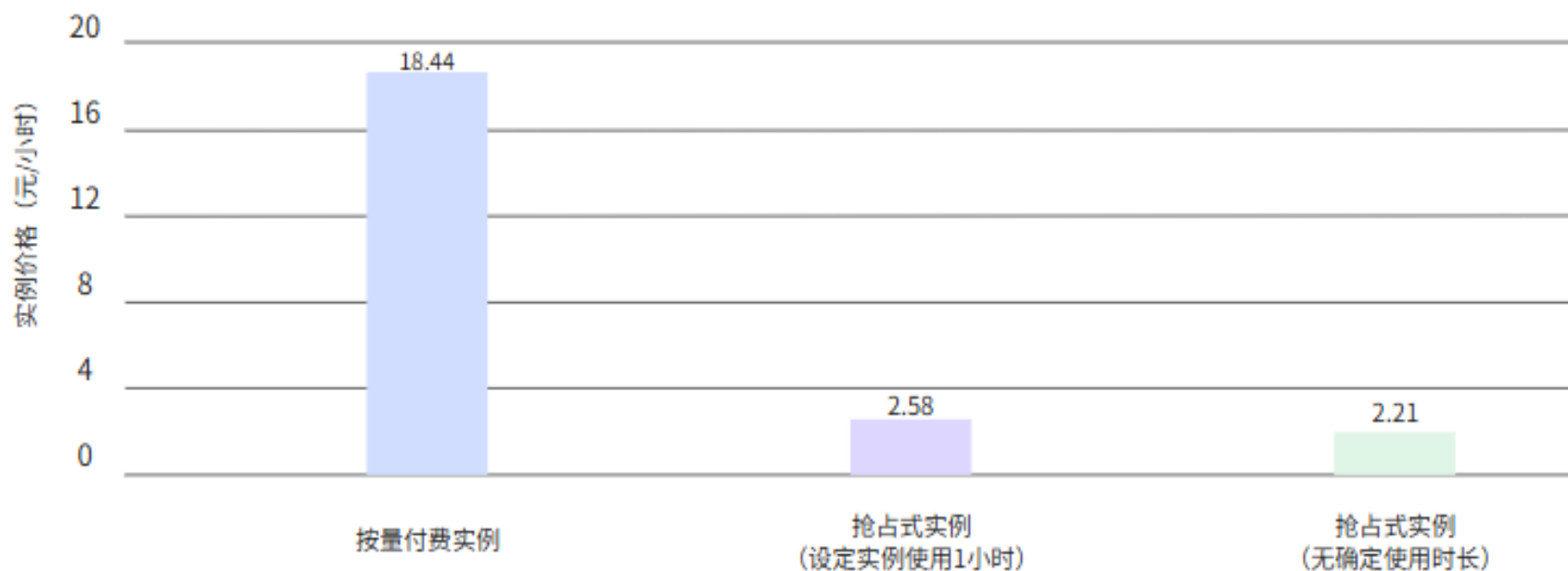
- high computational and memory requirements of LLMs make it challenging to serve them on hardware platforms

model	parameters	device	num	type
GPT-3	175 bilion	Nvidia A100-40GB	16	float

# Background : Spot Instances

Modern clouds offer:

- on-demand instances
- spot instances

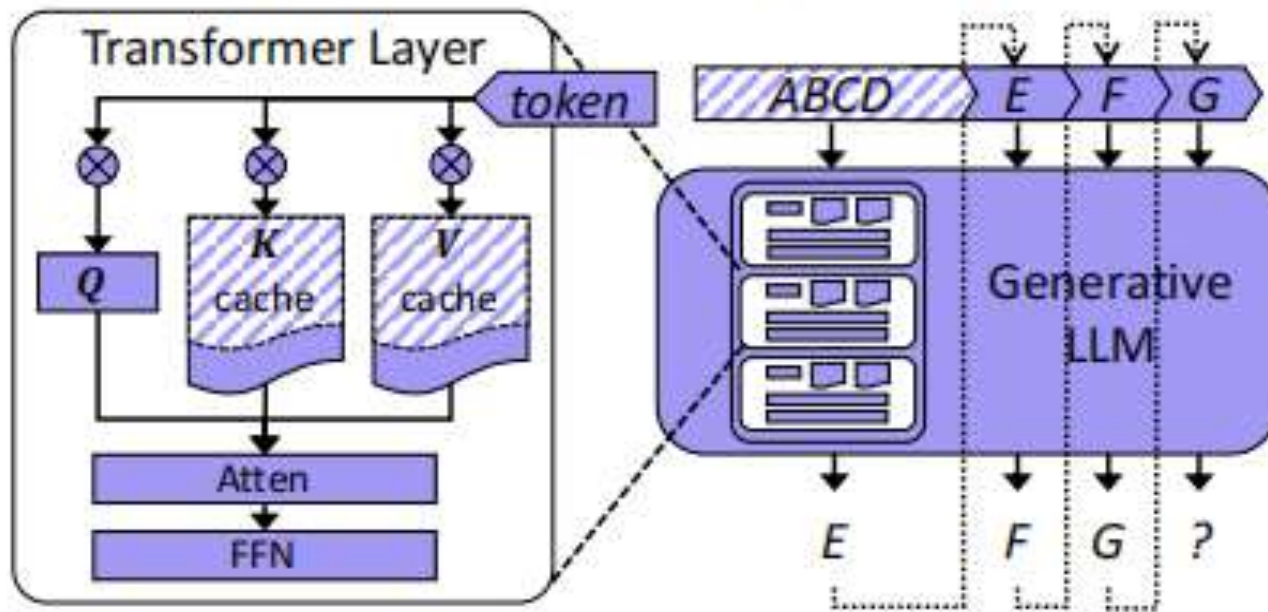


# Background: Generative LLM inference

- $l_{exe}$  is execution latency
- $t_{exe}$  is the LLM's execution time

Execution Latency:

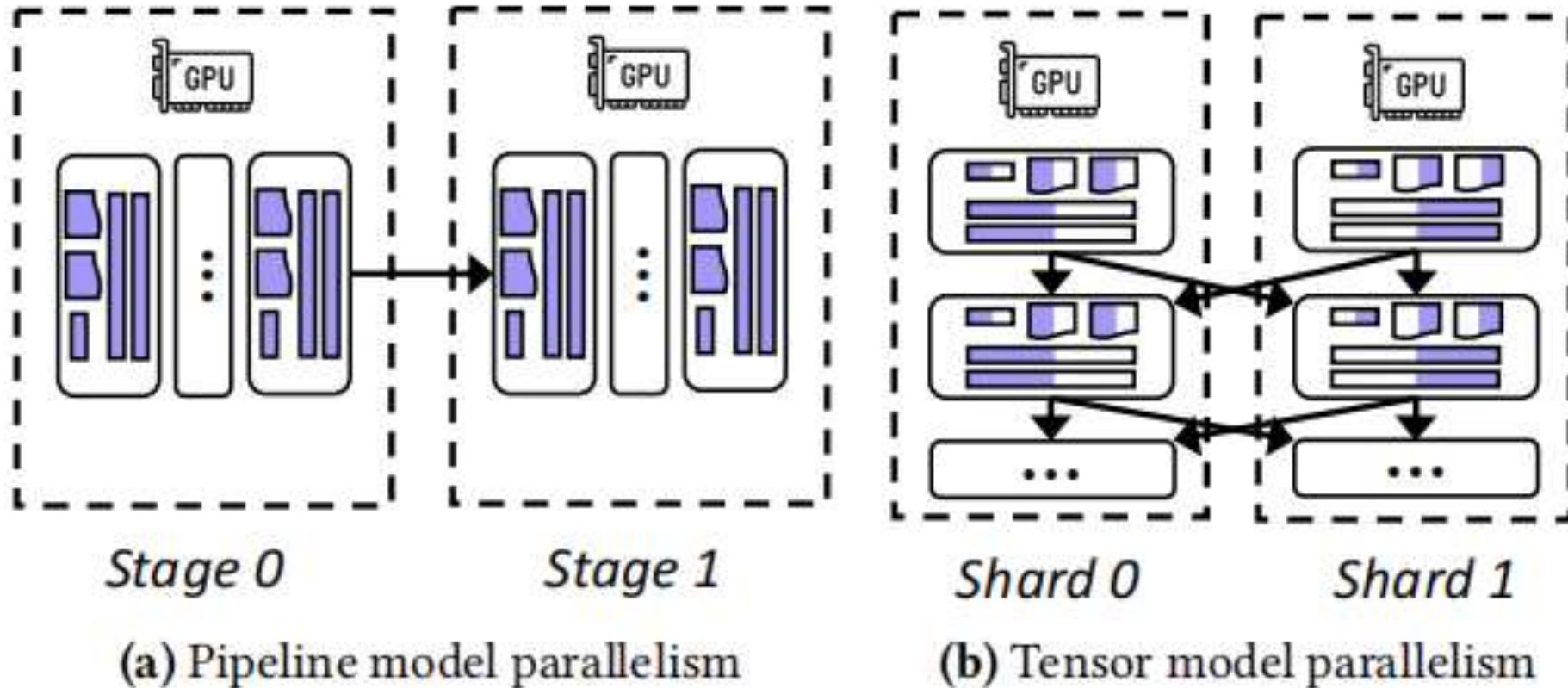
$$l_{exe}(3|4) = t_{exe}(4) + t_{exe}(1) + t_{exe}(1) + t_{exe}(1)$$



$$l_{exe}(S_{out}|S_{in}) = t_{exe}(S_{in}) + \sum_{i=1}^{S_{out}} t_{exe}(S_{in} + i)$$
$$\approx t_{exe}(S_{in}) + S_{out} \times t_{exe}(1)$$

# Background: Distributed Inference

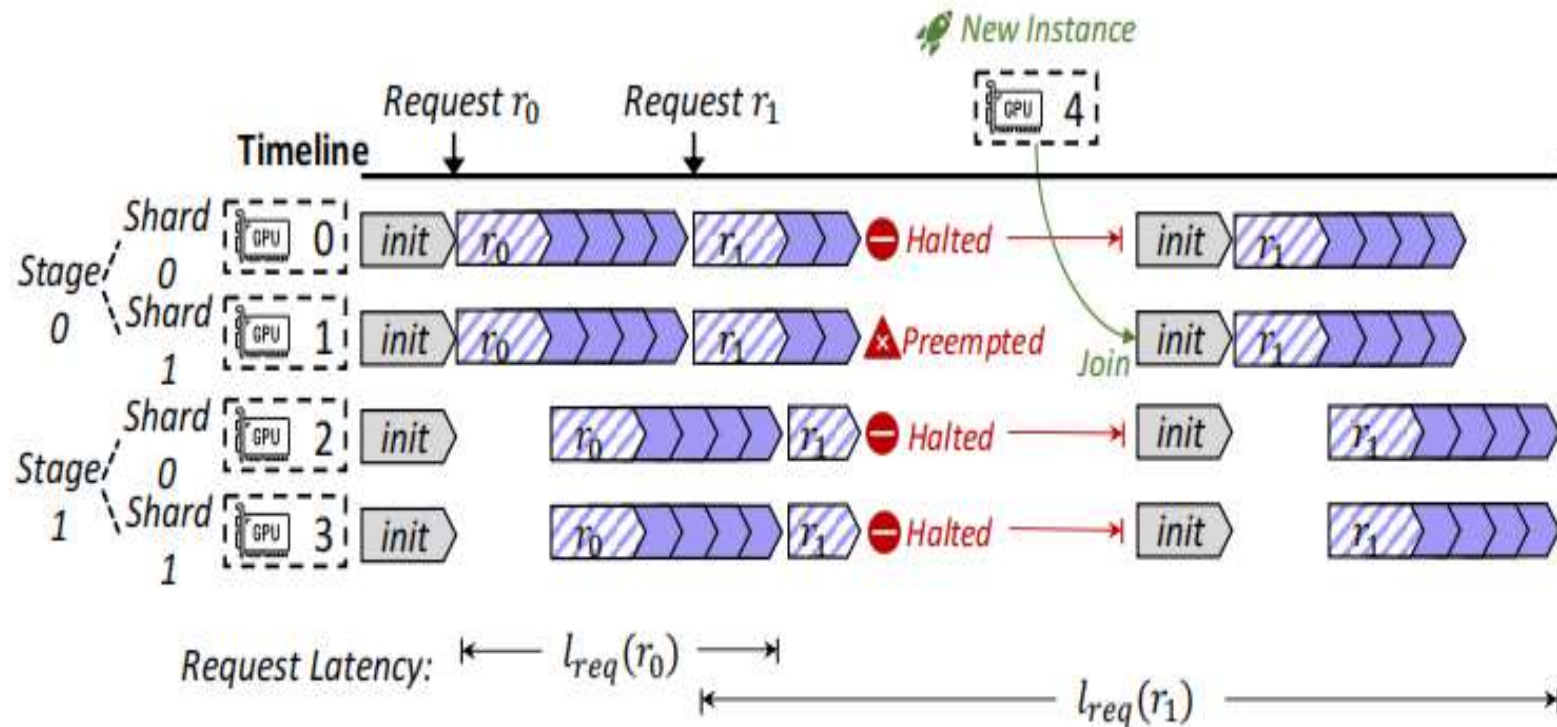
- Inference pipeline combine pipeline model parallelism and tensor model parallelism





# Background: Preemptible LLM Inference

- One Preempted hang all the other instances
- New instance joins need initialization costs



(b) Distributed LLM inference ( $P=2$ ,  $M=2$ , batch size=1) on preemptible instances



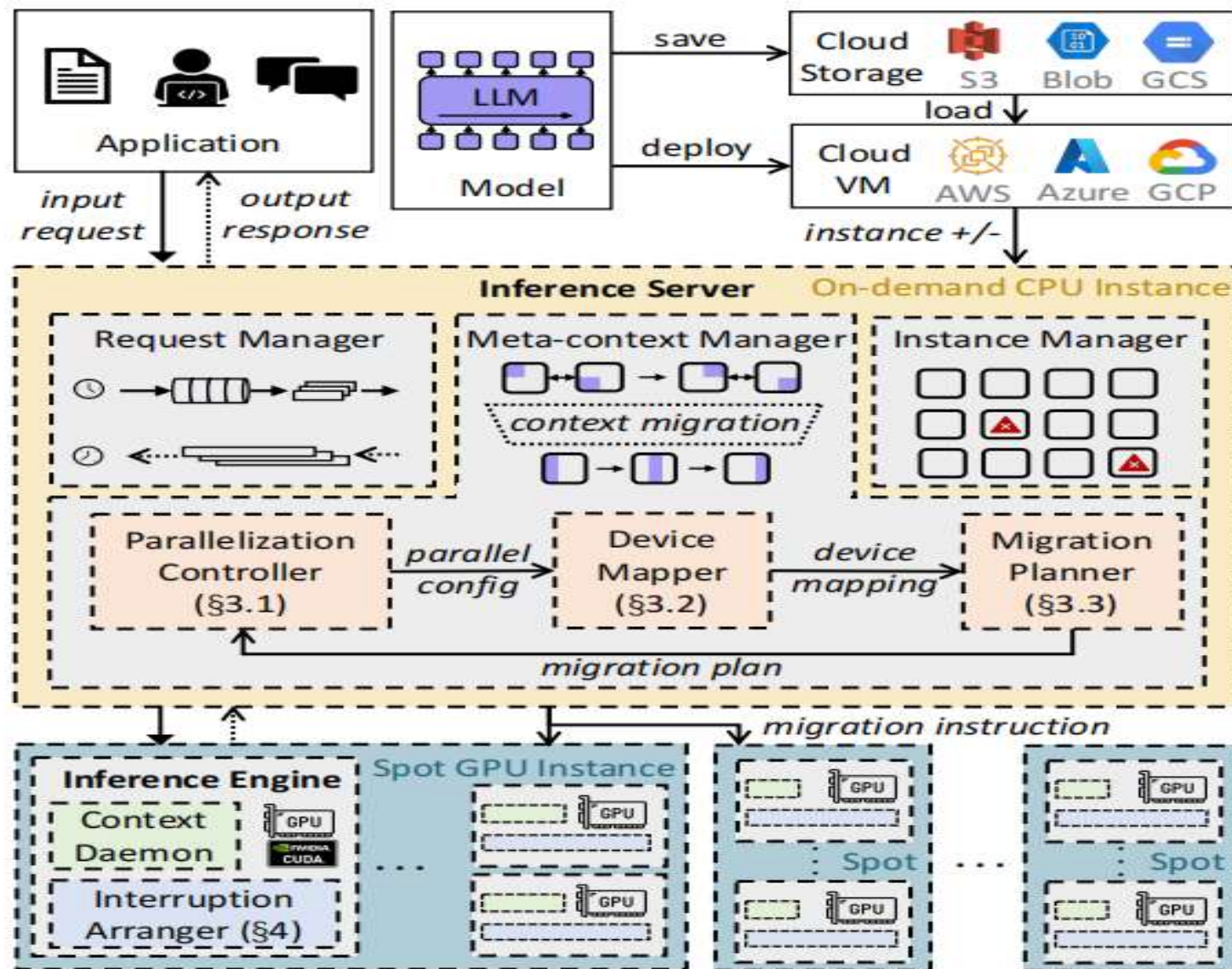
# Design:SpotServe

## Inference Server:

- Request Manager
- Meta-context Manager
- Instance Manager

## Inference Engine :

- Context Daemon
- Interruption Arranger



# Design: Parallelization Controller

- Parallelization Controller adjusts the parallelization configuration
- Algorithm 1 balance the trade-off among throughput latency, and cost

---

**Algorithm 1** Adaptive configuration optimizer.

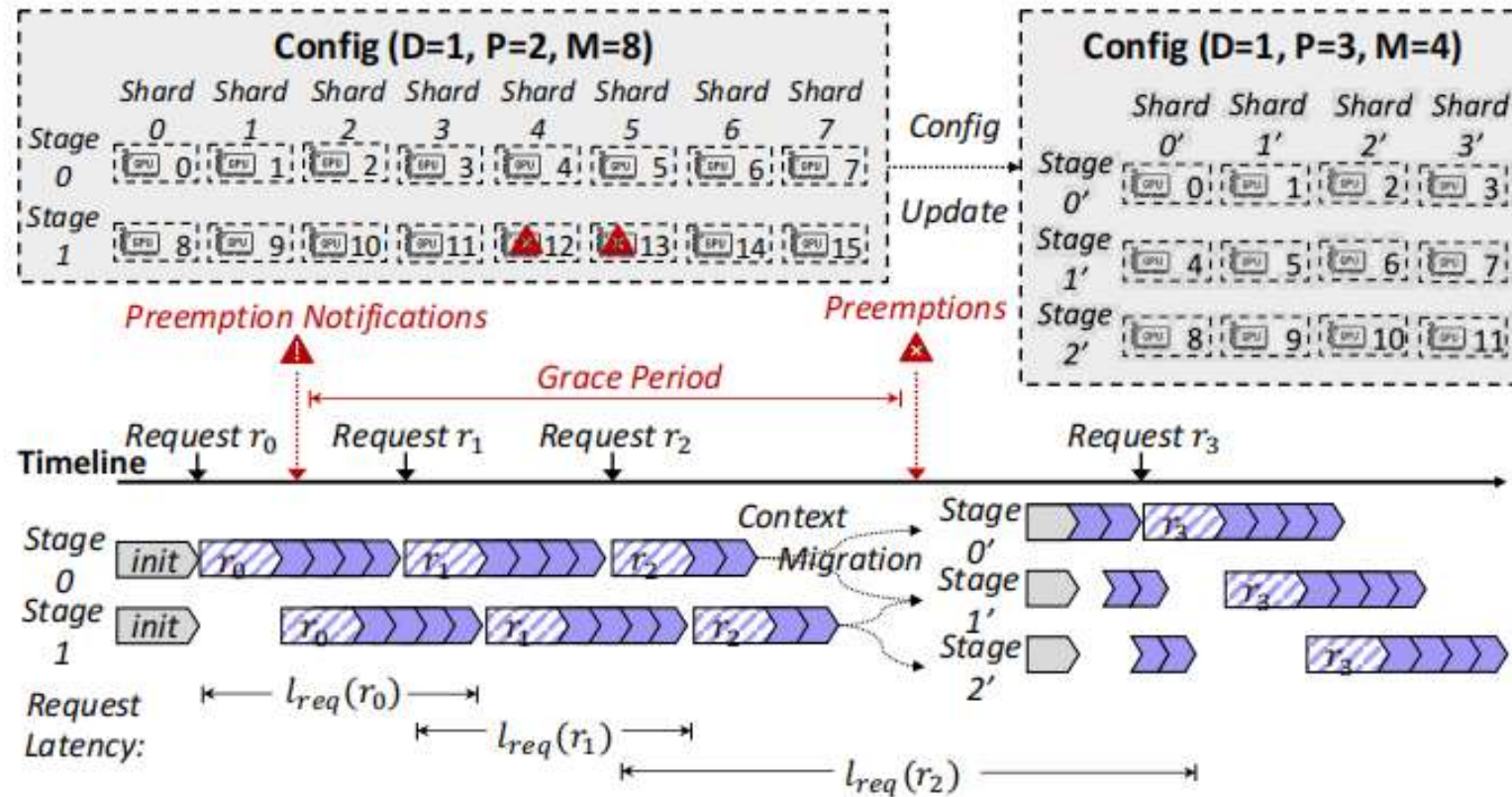
---

```
1: function CONFIGOPTIMIZER( $N_t, C_t, \alpha_t$ )
2:   if  $\exists C. \phi(C) \geq \alpha_t$  and cloud has enough instances for  $C$ 
   then
3:      $C_{t+1} \leftarrow \arg \min_{C | \phi(C) \geq \alpha_t} l_{req}(C)$ 
4:   else
5:      $C_{t+1} \leftarrow \arg \max_{C | N_t} \phi(C)$ 
6:    $\Delta \leftarrow \#Instances(C_{t+1}) - N_t$ 
7:   if  $\Delta > 0$  then
8:     InstanceManager.alloc( $\Delta$ , ondemand_and_spot)
9:   else
10:    InstanceManager.free( $-\Delta$ , ondemand_first)
11:  ConfigUpdate( $C_t, C_{t+1}$ )
```

---

# Design:Device Mapper

- A context migration mechanism can resume interrupted requests' inference

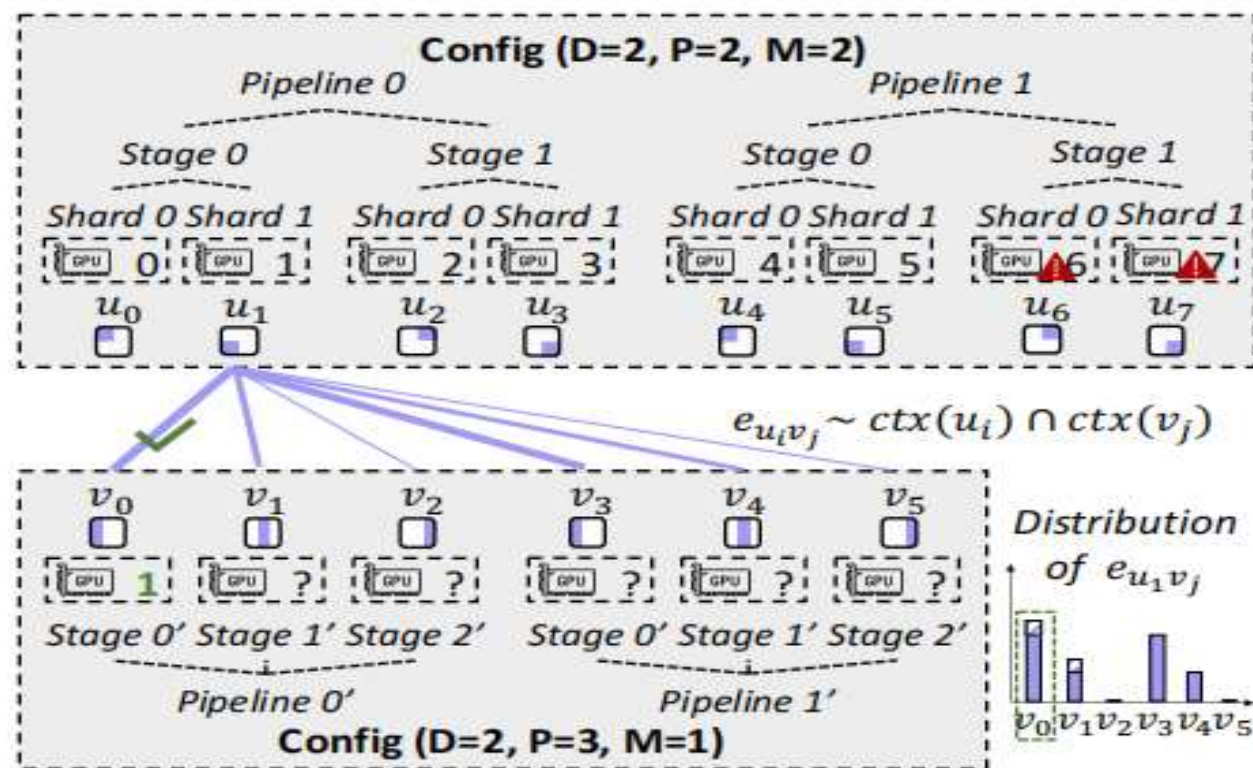


(a) Illustration of configuration update and context migration in SpotServe.



# Design:Device Mapper

- Device mapping problem  $\rightarrow$  bipartite graph matching task
- Maximally reuses the model parameters and KV cache and minimizes the total data transmission.



(b) Illustration of device mapping in SpotServe.

# Design:Migration Planner

- Algorithm2 is a progressive migration schedule
- Consider the memory usage during the progressive migration

---

**Algorithm 2** Workflow of the SpotServe migration planner.

---

```

  ▶ Progressive Migration
1: function MIGRATIONPLANNER(context ctx, plan = [ ])
2:   plan.append(<migrate, ctx.cache>)
3:    $O \leftarrow$  Layer migration order from MemOptMigPlanner
4:   for layer index  $i$  in range(0, #layers) do
5:     plan.append(<migrate, ctx.weight[ $O_i$ ]>)
6:      $p \leftarrow$  Get pipeline stage index of layer  $O_i$ 
7:     if stage  $p$  completes all context migration then
8:       plan.append(<start, instances of stage  $p$ >)
  ▶ Memory Optimized Migration
9: function MEMOPTMIGPLANNER(maximum buffer size  $U_{max}$ )
10:   $O \leftarrow [ ]$ ,  $X \leftarrow \{ \}$ 
11:  Instance buffer memory usage  $U \in \{0\}^N$ 
12:  for layer index  $i$  in range(0, #layers) do
13:    if (migrate, ctx.weight[ $i$ ]) doesn't exceed  $U_{max}$  then
14:      Update buffer memory usage  $U$ 
15:       $O.append(i)$ 
16:    else
17:       $X.add(i)$ 
18:  while  $X$  is not empty do
19:     $x_{opt} \leftarrow$ 
20:       $\arg \min_{x \in X} \max_{0 \leq i \leq N-1} \{U_i \mid (\text{migrate}, \text{ctx.weight}[x])\}$ 
21:     $O.append(x_{opt})$ 
     $X.remove(x_{opt})$ 
```

---

# Design:JIT Arrangement

Utilize grace period

- Each spot GPU instance includes an interruption arranger
- Just-in-time (JIT) arrangement decide when to stop decoding

$$S_t = \arg \max_{0 \leq s \leq s_{out}} \{l_{exe}(S|C_t) < T^- - T_{mig}\} \\ \cap T_{mig} < l_{exe}(S_t|C_t)$$

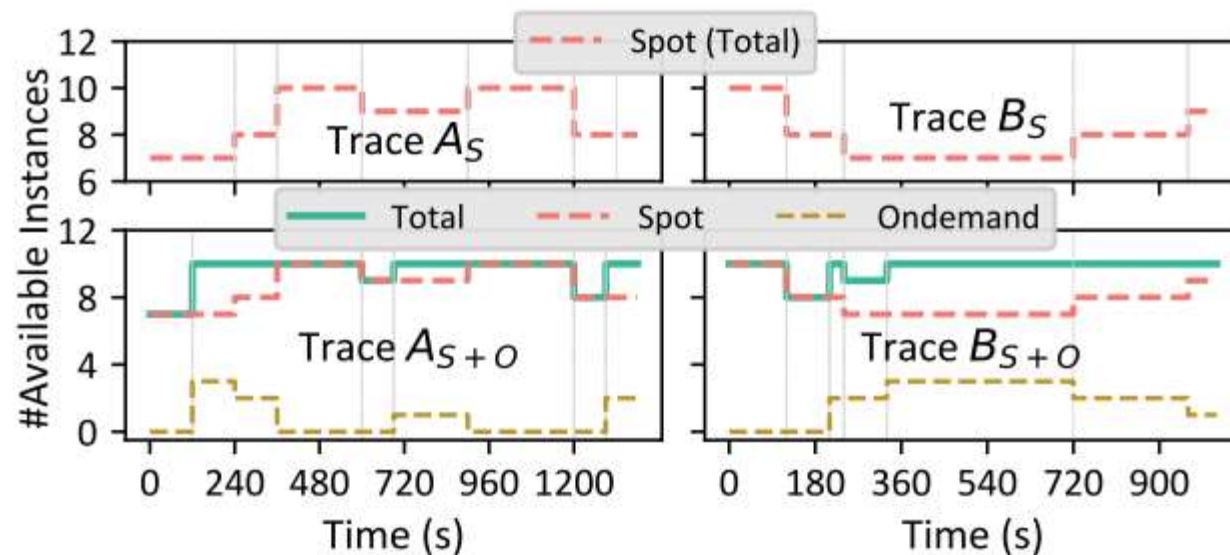
# Evaluation: Experiment Setup

## Baseline:

- Request rerouting
- Model reparallelization

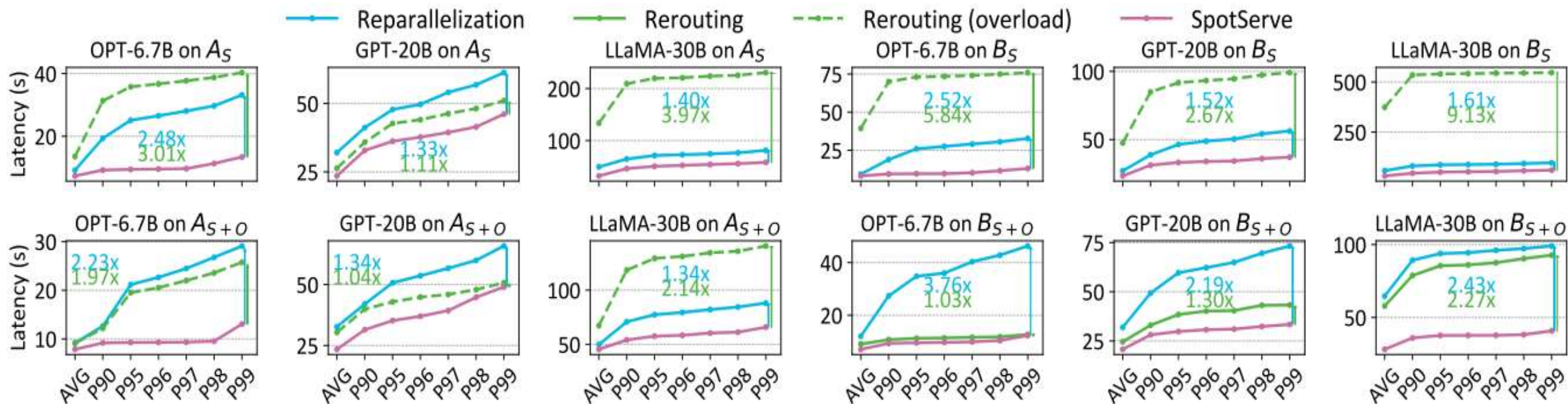
## Models:

- OPT-6.7B
- GPT-20B
- LLaMA-30B

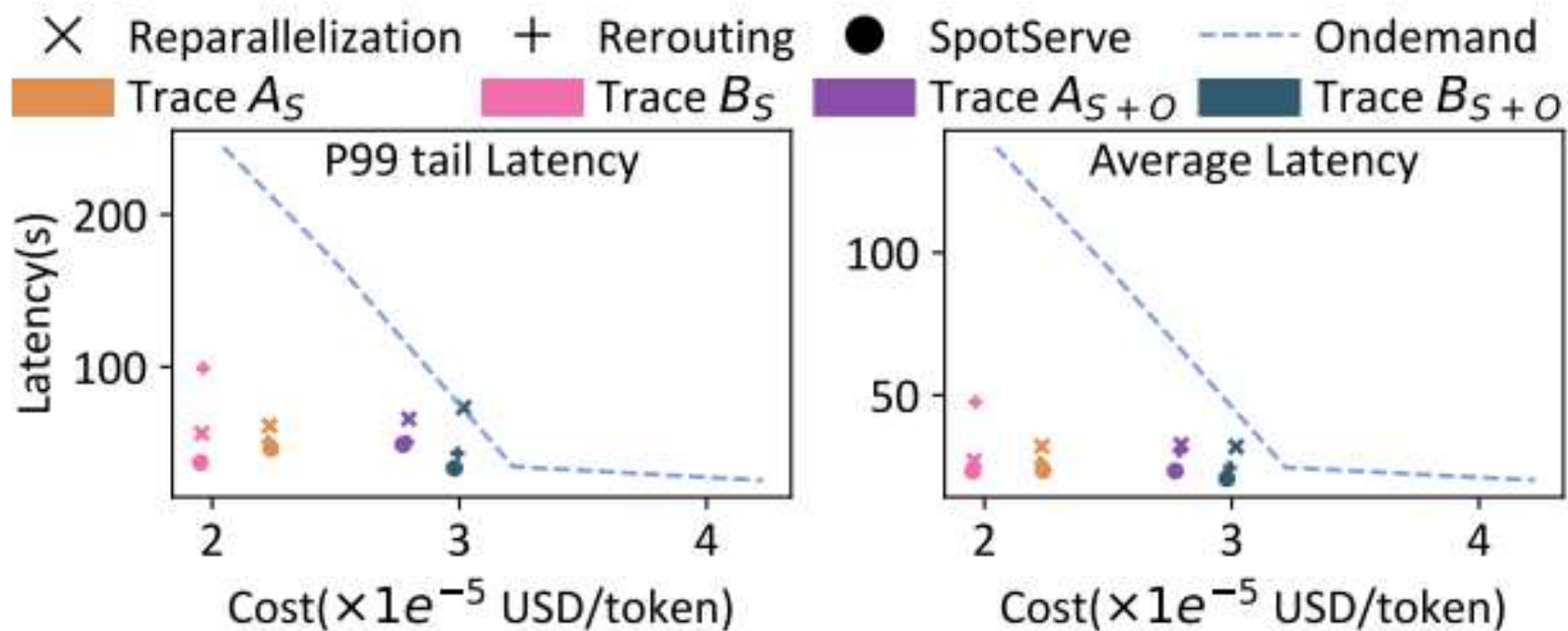




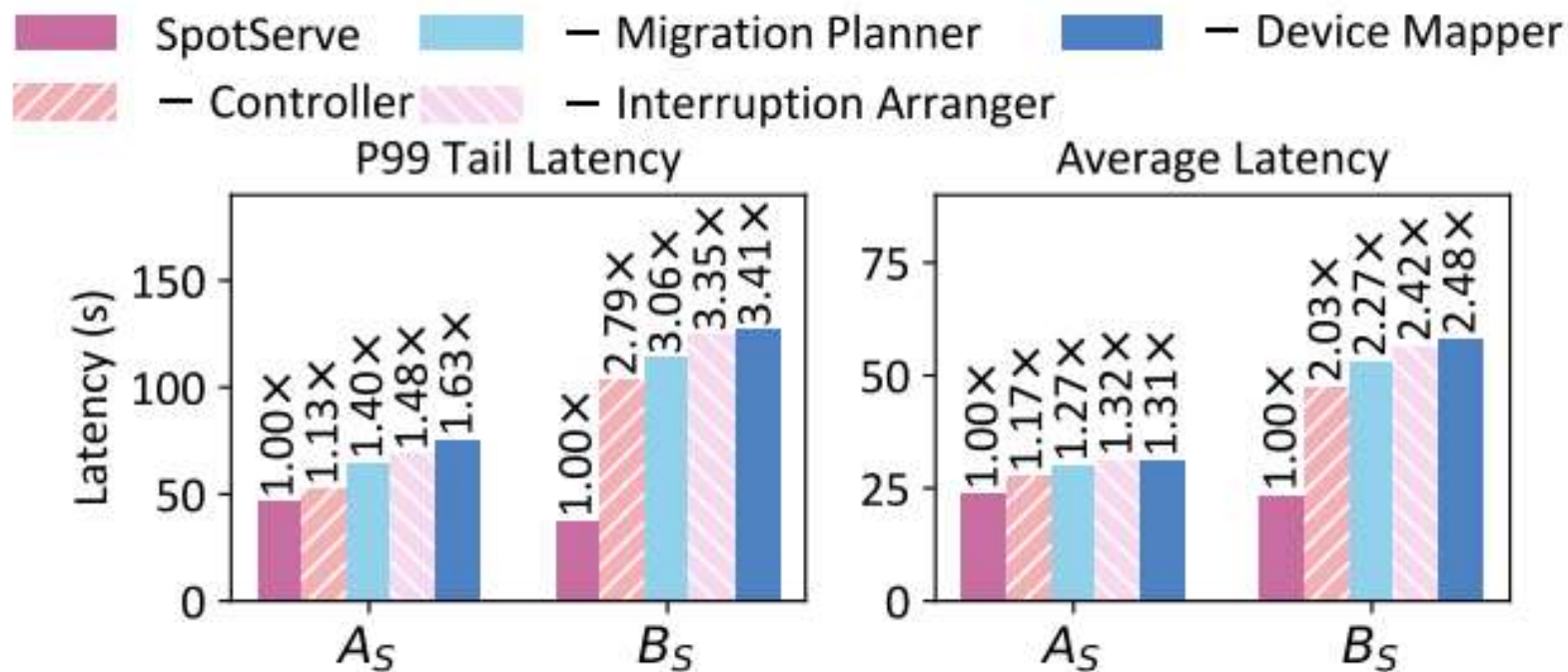
# Evaluation



# Evaluation



# Evaluation



# Thinking

- 文章有什么问题
  1. 该系统的并行策略中考虑的都是相同类型的GPU实例，并行策略上也是将其按照同构资源处理
  2. 在模型映射部分，只将模型参数的可复用程度作为边权的考虑，没有考虑设备间数据传输的能力差异
- 该方法能不能应用到别的场景

对于智能家居的分布式计算来说，对于一个有处理任务的设备来说，其他当前状态下空闲的设备的资源对它来说视作可抢占资源，利用目前可用资源配置出最优的并行推理流水线

Thanks