# LLMServingSim: A Simulation Infrastructure for LLM Inference Serving Systems

**Jaehong Cho**

Minsu Kim

Hyunmin Choi

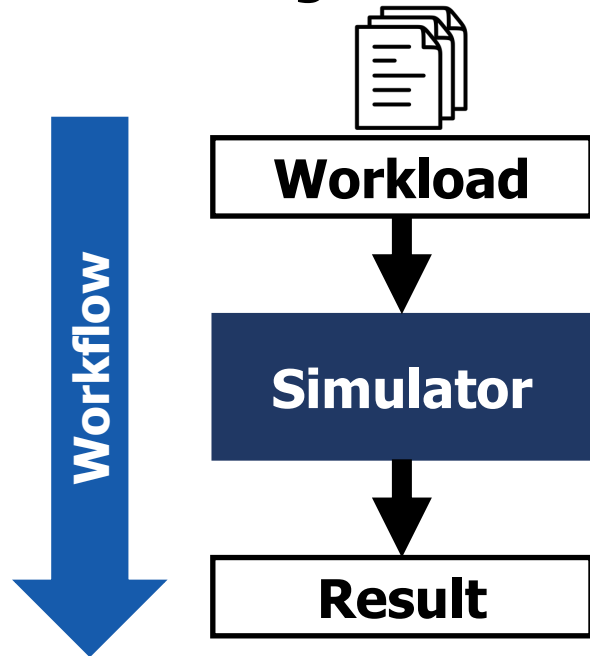Jongse Park

# Large Language Model (LLM)
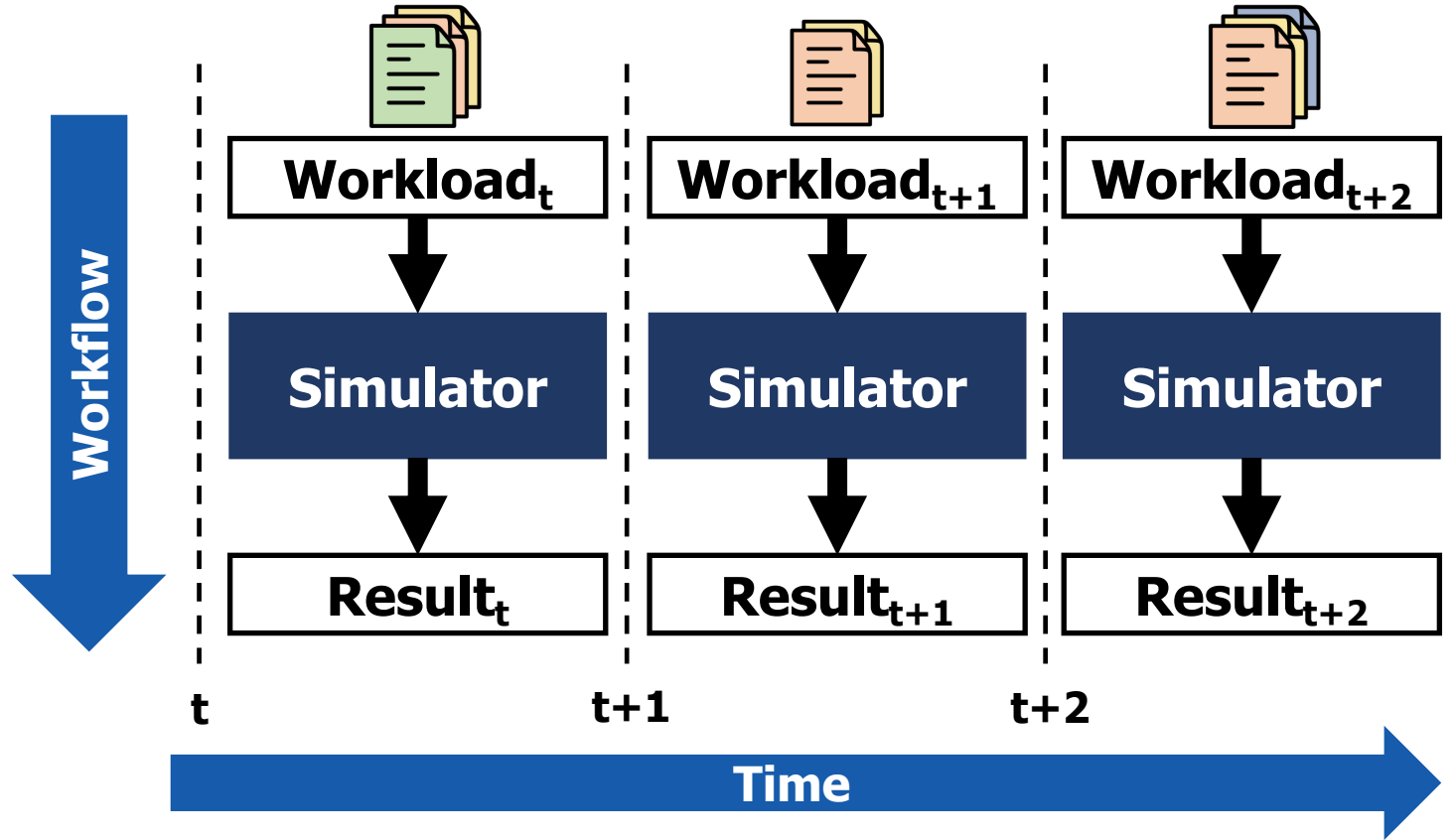
## LLM Inference Serving System

**Input Prompt**

"Large Language Model"

**LLM Inference Serving System**

**Output**

"Large Language Model is awesome and used in many real-world applications."

# Challenge 1: Autoregressive LLM

- **Repetitive identical iterations for "training"**

- **Dynamic workloads differ in time**



**Workflow**

Workload → Simulator → Result

**Workflow**

$Workload_t$ → Simulator → $Result_t$

$Workload_{t+1}$ → Simulator → $Result_{t+1}$

$Workload_{t+2}$ → Simulator → $Result_{t+2}$

t          t+1          t+2

**Time**

# LLM Inference Serving
## Iteration-level Scheduling & KV Cache Paging

Block: Unit of managing KV cache

**Execution Engine**

**Batch**

**Iteration-level**

**Result**

**Request Pool**

**New Request**

**End Request**

**Logical KV cache blocks**

| | | | | |
|---|---|---|---|---|
| Block 0 | Large | Language | Model | is |
| Block 1 | awesome | and | used | in |
| Block 2 | many | | | |

**Block Table**

**GPU Physical KV cache blocks**

| | | | | | |
|---|---|---|---|---|---|
| Block 0 | Large | Language | Model | is | Allocation |
| Block 1 | many | | | | |
| Block 2 | awesome | and | used | in | |

**CPU Physical KV cache blocks**

| | | | | | |
|---|---|---|---|---|---|
| Block 0 | Computer | architecture | and | system | Eviction |

KAIST School of Computing

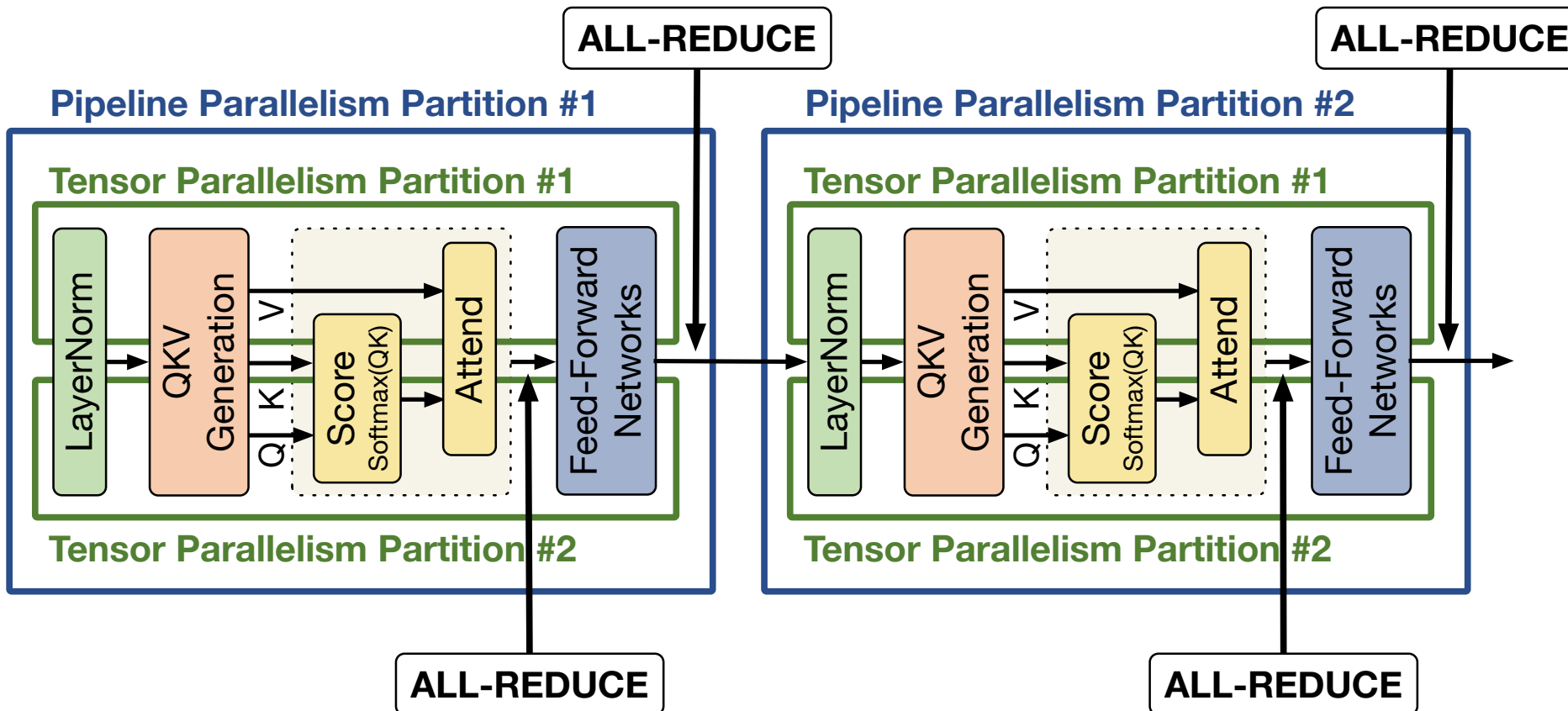CASYS | KAIST Computer Architecture & System Lab

# Solution 1: Iteration-level Scheduling

# Solution 1: KV Cache Paging

# Challenge 2: LLM Specific Parallelism

# Solution 2: LLM Graph Converter
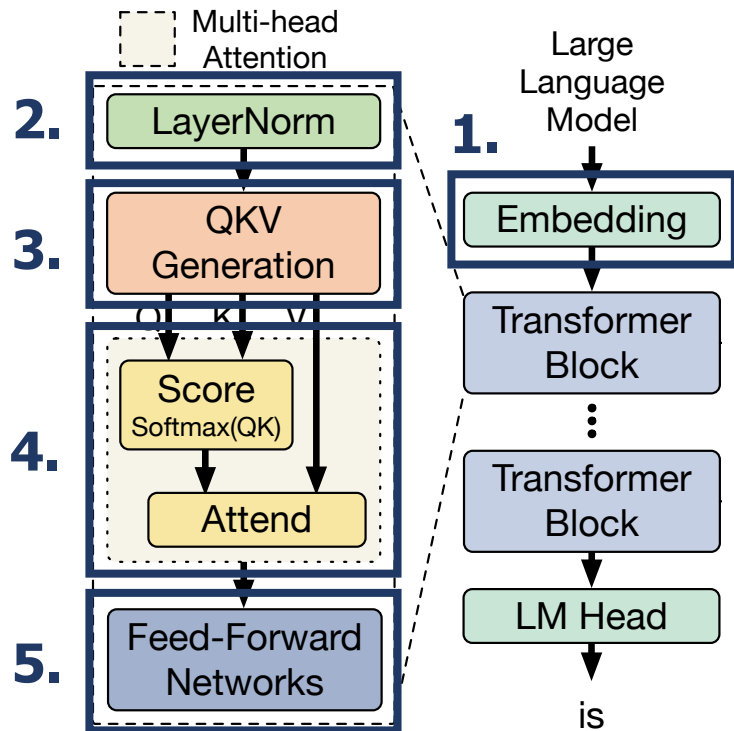


**Pipeline Parallelism**

**Tensor Parallelism**

Host CPU

| ID | Op | Src | Dest |
|----|------|---------|---------|
| 0 | Embd | - | - |
| 1 | | - | - |
| ... | | - | - |
| N | **Send** | Group 1 | Group 2 |

| ID | Op | Src | Dest |
|----|------|---------|---------|
| 0 | **Recv** | Group 1 | Group 2 |
| 1 | | - | - |
| ... | | - | - |
| N | **Send** | Group 2 | Group 3 |

| ID | Op | Src | Dest |
|----|------------|-----|------|
| ... | | - | - |
| k | **ALL-REDUCE** | ALL | ALL |
| ... | | - | - |

# Challenge 3: Slow Simulation Time

- Batch 32

- Sequence Length 512

- 1 Iteration



More than **5 hours** for 1 iteration

More than **2 hours** for 1 iteration

# Solution 3: Computation Reuse



Multi-head Attention

2. LayerNorm

3. QKV Generation

Q  K  V

4. Score Softmax(QK) → Attend

5. Feed-Forward Networks

1. Large Language Model

Embedding

Transformer Block

⋮

Transformer Block

LM Head

is

**Make Traces of Each Layer**

1. 2. 3. 4. 5.

1. [2. 3. 4. 5.] ... [2. 3. 4. 5.]

Number of Transformer Blocks

**Make Full Model Trace**

Cache

Load

1. [2. 3. 4. 5.] ... [2. 3. 4. 5.]

**Swap Attention Layer & Reuse**
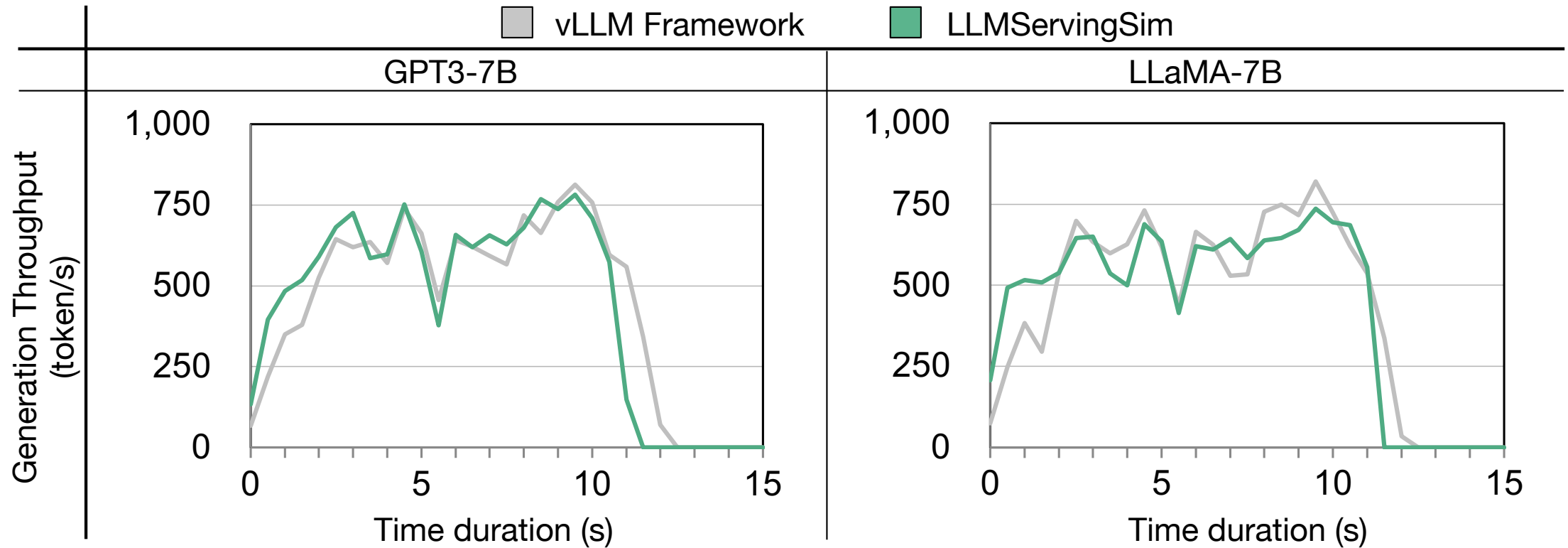
KAIST School of Computing

CASYS | KAIST Computer Architecture & System Lab
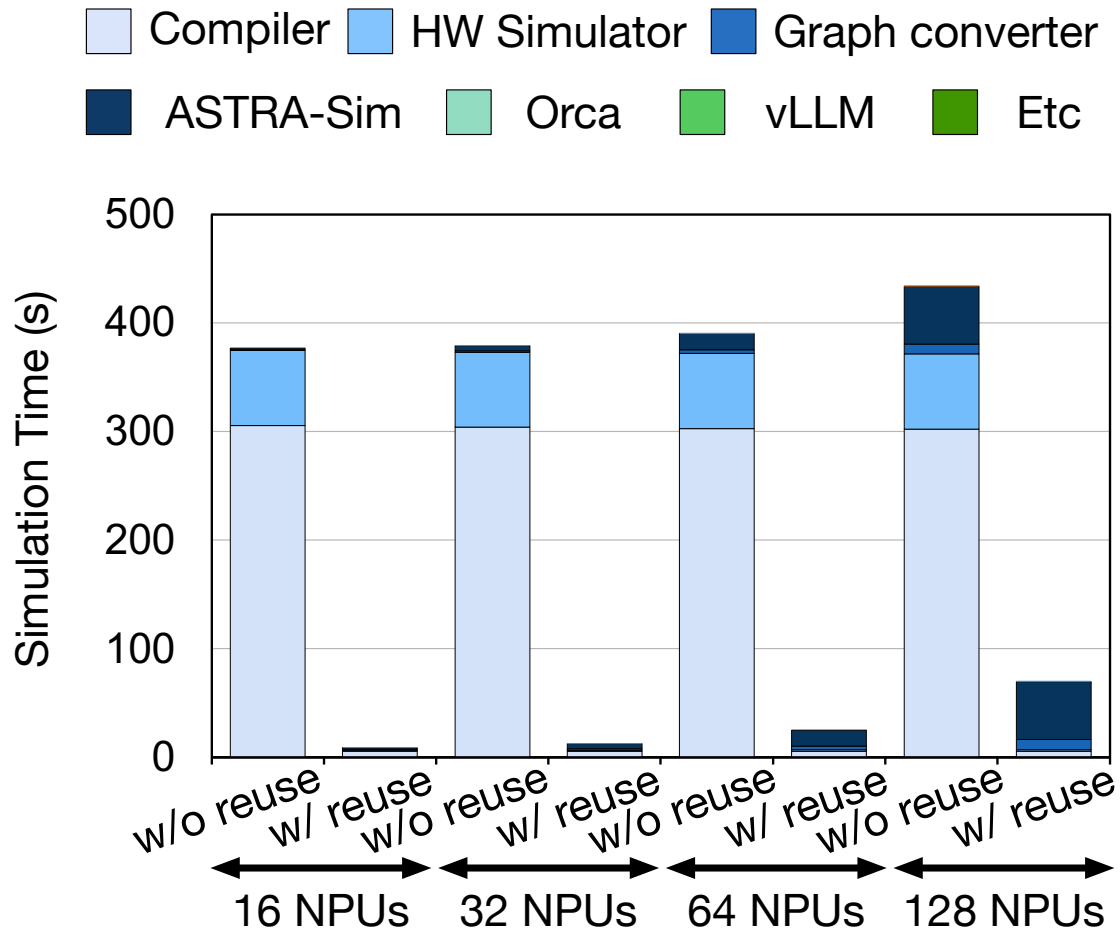
# Evaluation

## Validation



- ShareGPT sampled requests, request arrival pattern using Poisson distribution
- Average error rate **14.7%**

# Evaluation

## Performance of Computation Reuse

Legend:
- Compiler
- HW Simulator
- Graph converter
- ASTRA-Sim
- Orca
- vLLM
- Etc



- GPT-3 175B 1 iteration
- Sequence Length: 2048

- Without reuse: **~400 sec**
- With reuse: **<1 min**

- **18.7×** speedup using computation reuse

# Conclusion

- **Large scale LLM inference serving system simulator**

- **14.7%** error rate against real GPU-based LLM serving system
- **18.7x** speedup using computation reuse