# NeuPIMs: NPU-PIM Heterogeneous Acceleration for Batched LLM Inferencing

**Guseul Heo**
Sangyeop Lee
Jaehong Cho
Hyunmin Choi
Sanghyeon Lee
Hyungkyu Ham[†]
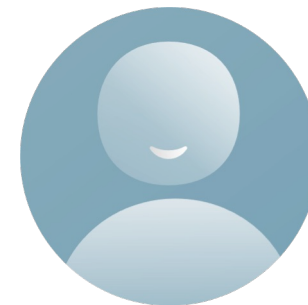Gwangsun Kim[†]
Divya Mahajan[§]
Jongse Park

KAIST

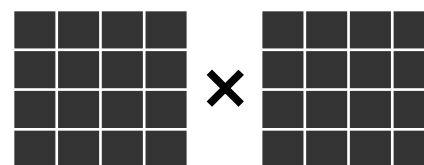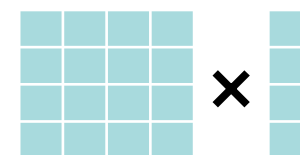POSTECH[†]

Georgia Institute of Technology[§]

KAIST

CASYS

ASPLOS 2024

LLM batched inference comprises **GEMM** and **GEMV**
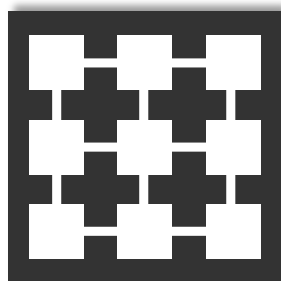
GEMM

matrix-matrix
multiplication

GEMV

matrix-vector
multiplication

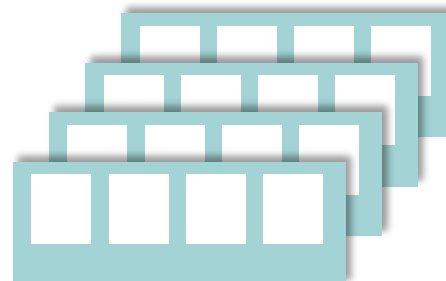However, with naïve NPU+PIM integration,
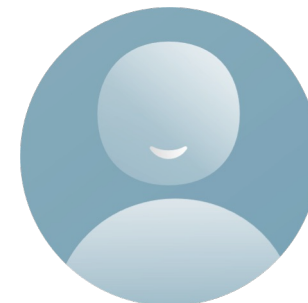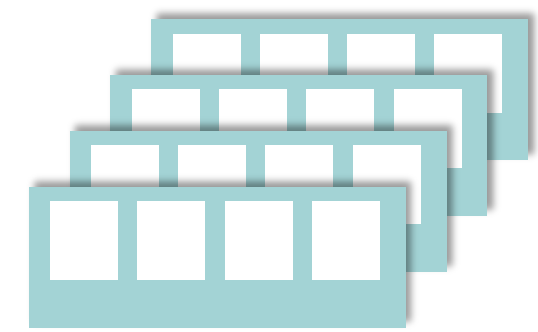
system suffers from resource underutilization

NPU + PIM

# Challenge #2

GEMM and GEMV have algorithmic dependency in LLM
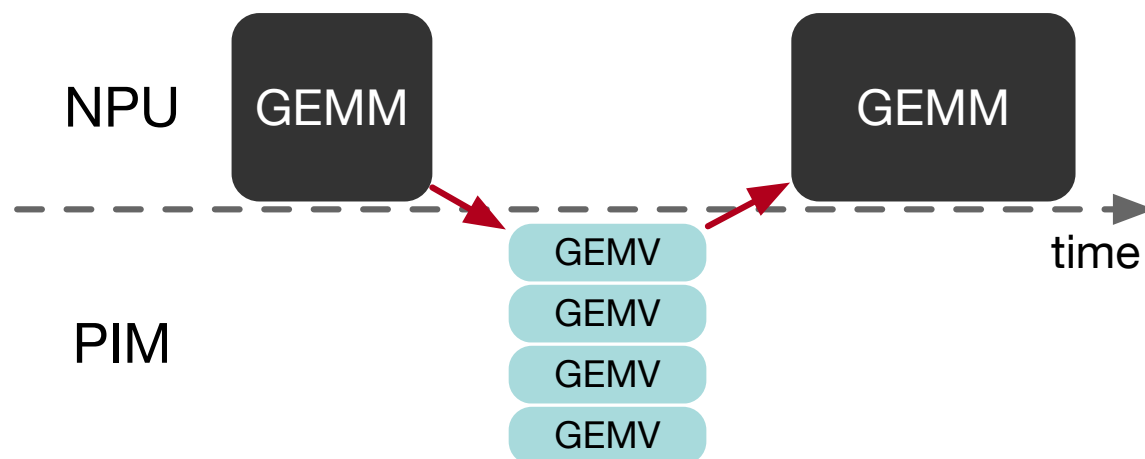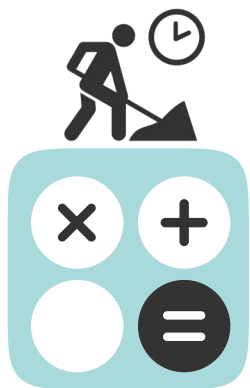
# We devise **NeuPIMs**, NPU-PIM heterogeneous acceleration solution for batched LLM inferencing
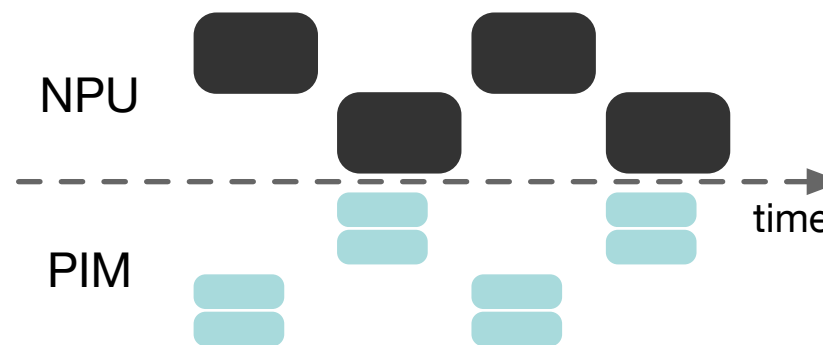
PIM

MEM

**PIM with dual row buffers**

NPU

PIM

time

**Sub-batch interleaving**

# NeuPIMs: NPU-PIM Heterogeneous Acceleration for Batched LLM Inferencing

**Session 6B**
04/30 (Tue)
14:30

**2.4×**
throughput improvement over NPU

**1.6×**
throughput improvement over naïve NPU+PIM

ASPLOS 2024

Our simulator code is available
https://github.com/casys-kaist/NeuPIMs