# PIMMiner: A High-performance PIM Architecture-aware Graph Mining Framework

Presenter: Jiya Su     Advisor: Rujia Wang

Computer Science of Illinois Institute of Technology

MICRO-55
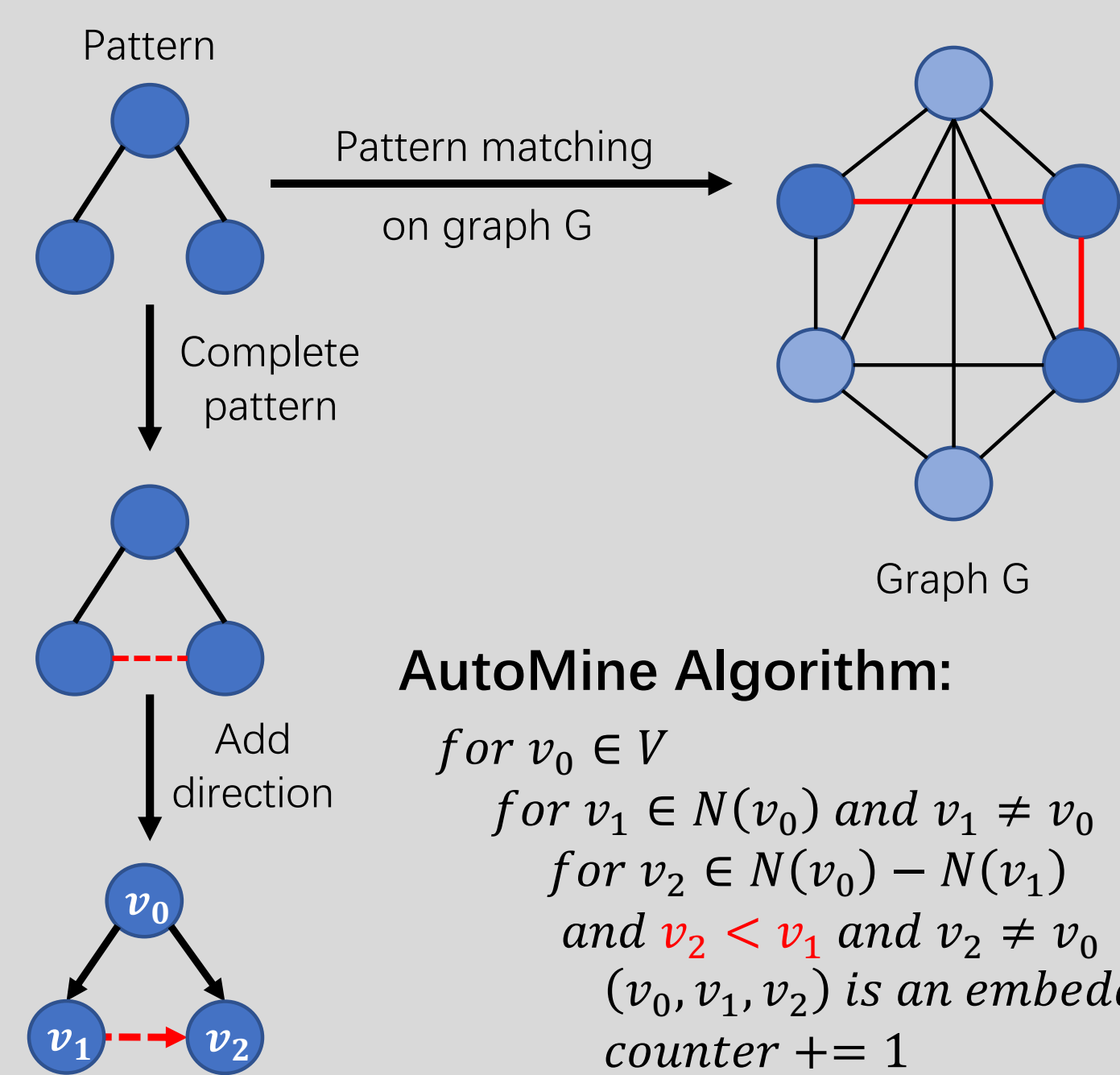
ILLINOIS INSTITUTE OF TECHNOLOGY

## Abstract

We propose PIMMiner, a high-performance PIM architecture graph mining framework.
- We first identify that current PIM architecture cannot be fully utilized by graph mining applications.
- Next, we propose a set of optimizations that enhance the locality, and internal bandwidth utilization and reduce remote bank accesses and load imbalance through cohesive algorithm and architecture co-designs.
- We compare PIMMiner with several state-of-the-art graph mining frameworks and show that PIMMiner is able to outperform all of them significantly.
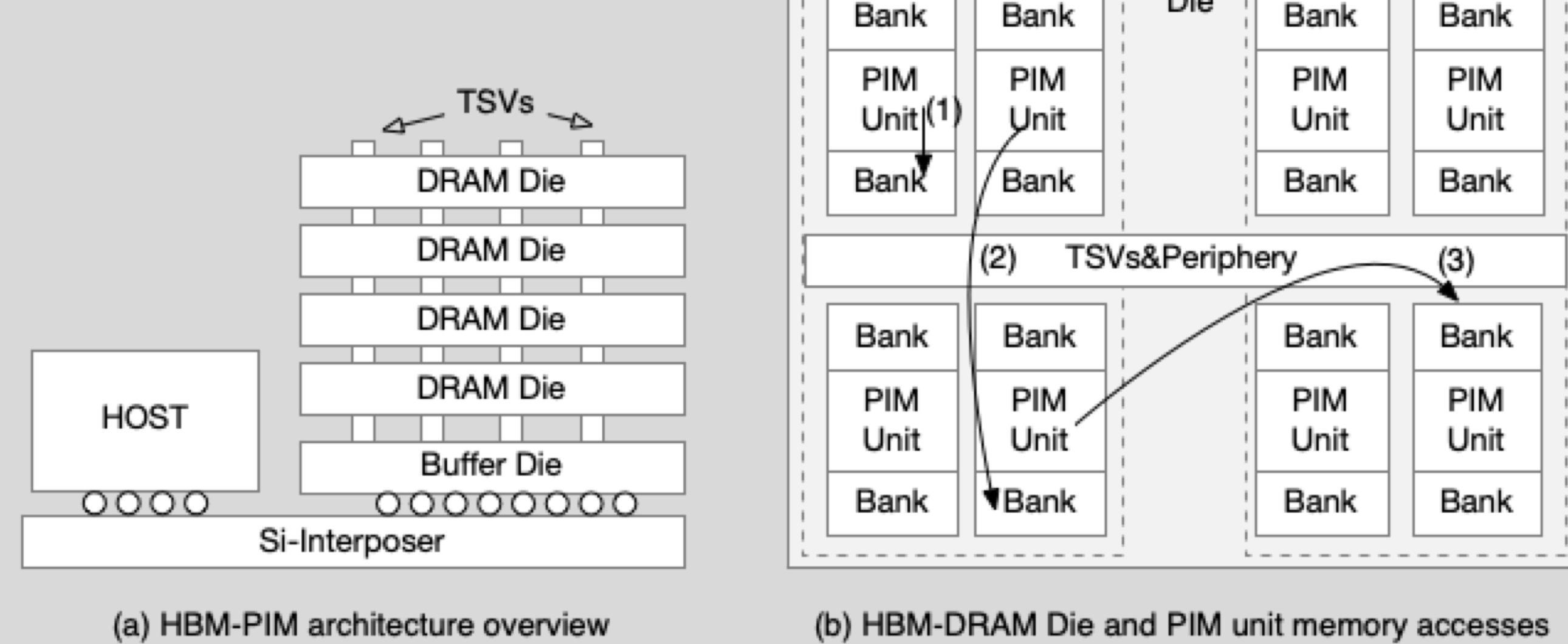
## Background

- Graph pattern mining (GPMI) needs to find all subgraphs with different patterns that meet the application requirements.
- GPMI applications are considered as a new class of data-intensive applications that generate massive irregular computation workloads and memory accesses, which degrade the performance significantly.

Pattern

Pattern matching on graph G

Complete pattern

Add direction

Graph G

AutoMine Algorithm:

$for\ v_0 \in V$
$\quad for\ v_1 \in N(v_0)\ and\ v_1 \neq v_0$
$\quad\quad for\ v_2 \in N(v_0) - N(v_1)$
$\quad\quad\quad and\ v_2 < v_1\ and\ v_2 \neq v_0$
$\quad\quad\quad (v_0, v_1, v_2)\ is\ an\ embedding$
$\quad\quad\quad counter\ += 1$

- Processing-in-Memory (PIM) integrates processing units inside the memory to reduce the overhead of frequent data movement and achieve high-performance and energy-efficient computation.
- Samsung has recently started manufacturing HBM-PIM chips. The HBM-PIM incorporates PIM cores inside of memory banks. There are three ways for a PIM unit to access memory: (1) near-core bank access, (2) intra-channel bank access; (3) inter-channel remote bank access.
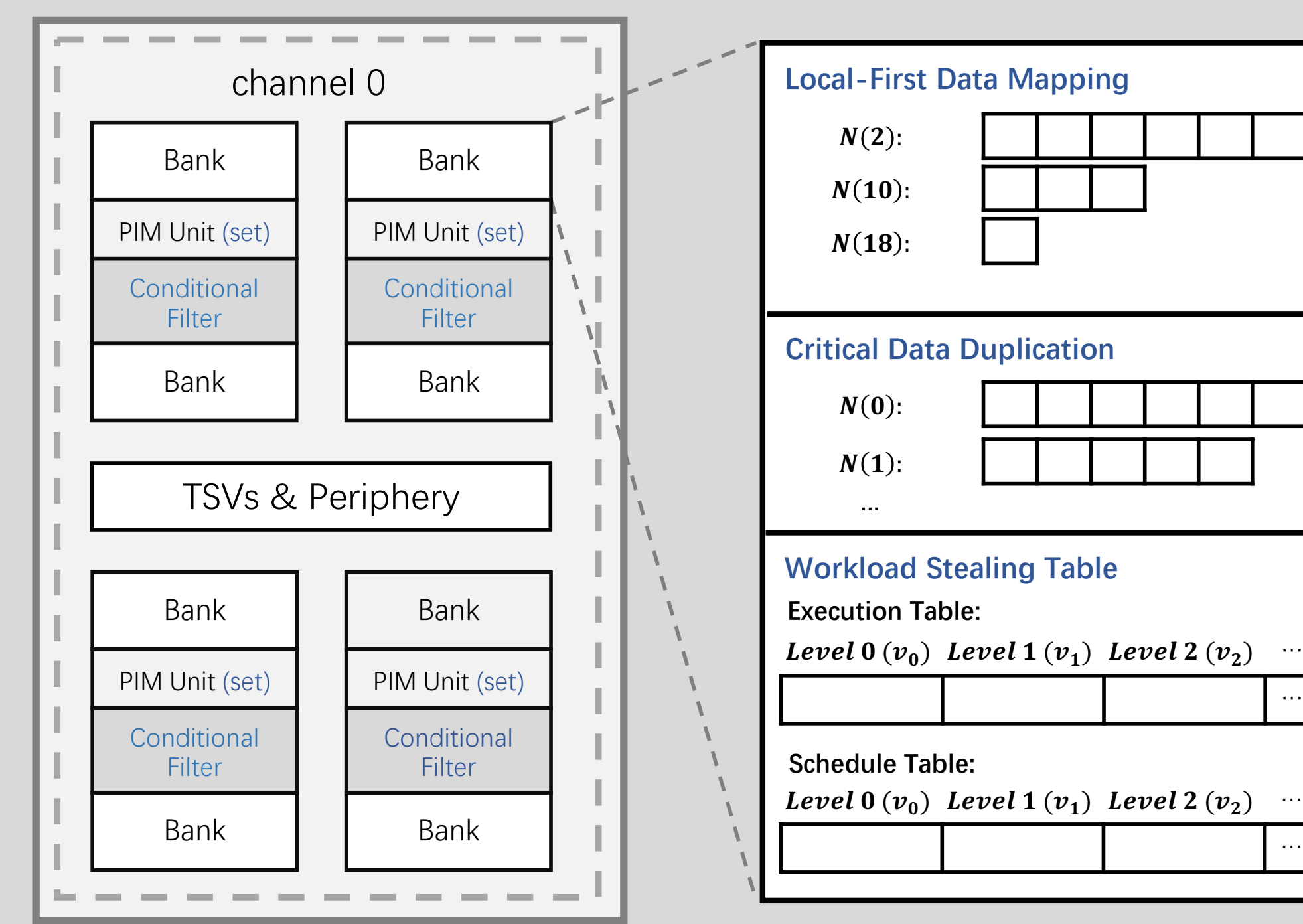
## Motivation

- Directly offload the GPMI execution kernel to PIM cannot achieve desired performance. We observe **high load imbalance** and lots of **inter-channel remote bank accesses**.

TABLE I: Performance between 96-threads CPU (L3 cache) and 128-core PIM (L1 cache) in 4-CC.

| Graph | Processing Unite | Execution Time (s) | Average Time (s) | Exe/Avg | Execution Speedup | Average Speedup |
|---|---|---|---|---|---|---|
| CI | CPU | 3.80E-05 | 1.65E-05 | 2.30 | 0.92x | 0.46x |
| | PIM | 4.13E-05 | 3.57E-05 | 1.16 | | |
| PP | CPU | 7.07E-04 | 3.37E-04 | 2.10 | 4.15x | 2.16x |
| | PIM | 1.71E-04 | 1.56E-04 | 1.09 | | |
| AS | CPU | 1.04E-02 | 7.04E-03 | 1.48 | 2.58x | 2.92x |
| | PIM | 4.05E-03 | 2.41E-03 | 1.68 | | |
| MI | CPU | 1.64E-01 | 8.48E-02 | 1.93 | 0.39x | 0.51x |
| | PIM | 4.22E-01 | 1.66E-01 | 2.53 | | |
| YT | CPU | 1.49E-01 | 9.18E-02 | 1.63 | 0.08x | 0.24x |
| | PIM | 1.83 | 3.78E-01 | 4.84 | | |
| PA | CPU | 1.83E-01 | 1.37E-01 | 1.34 | 2.66x | 2.53x |
| | PIM | 6.89E-02 | 5.40E-02 | 1.28 | | |
| LJ | CPU | 3.09 | 2.59 | 1.19 | 0.03x | 0.56x |
| | PIM | 103.13 | 4.63 | 22.29 | | |

## PIMMiner Framework Overview

PIMMiner has five lightweight architectural optimizations that work cohesively.

channel 0

Bank / Bank
PIM Unit (set) / PIM Unit (set)
Conditional Filter / Conditional Filter
Bank / Bank

TSVs & Periphery

Bank / Bank
PIM Unit (set) / PIM Unit (set)
Conditional Filter / Conditional Filter
Bank / Bank

Local-First Data Mapping
$N(2)$:
$N(10)$:
$N(18)$:

Critical Data Duplication
$N(0)$:
$N(1)$:
...

Workload Stealing Table
Execution Table:
$Level\ 0\ (v_0)$  $Level\ 1\ (v_1)$  $Level\ 2\ (v_2)$ ...

Schedule Table:
$Level\ 0\ (v_0)$  $Level\ 1\ (v_1)$  $Level\ 2\ (v_2)$ ...

## HBM-PIM Architecture

TSVs

DRAM Die
DRAM Die
DRAM Die
DRAM Die
Buffer Die
Si-Interposer

HOST

channel 0
Bank / Bank
PIM Unit (1) / PIM Unit
Bank / Bank

DRAM Die

channel 1
Bank / Bank
PIM Unit / PIM Unit
Bank / Bank

(2) TSVs&Periphery (3)

Bank / Bank
PIM Unit / PIM Unit
Bank / Bank

Bank / Bank
PIM Unit / PIM Unit
Bank / Bank

(a) HBM-PIM architecture overview
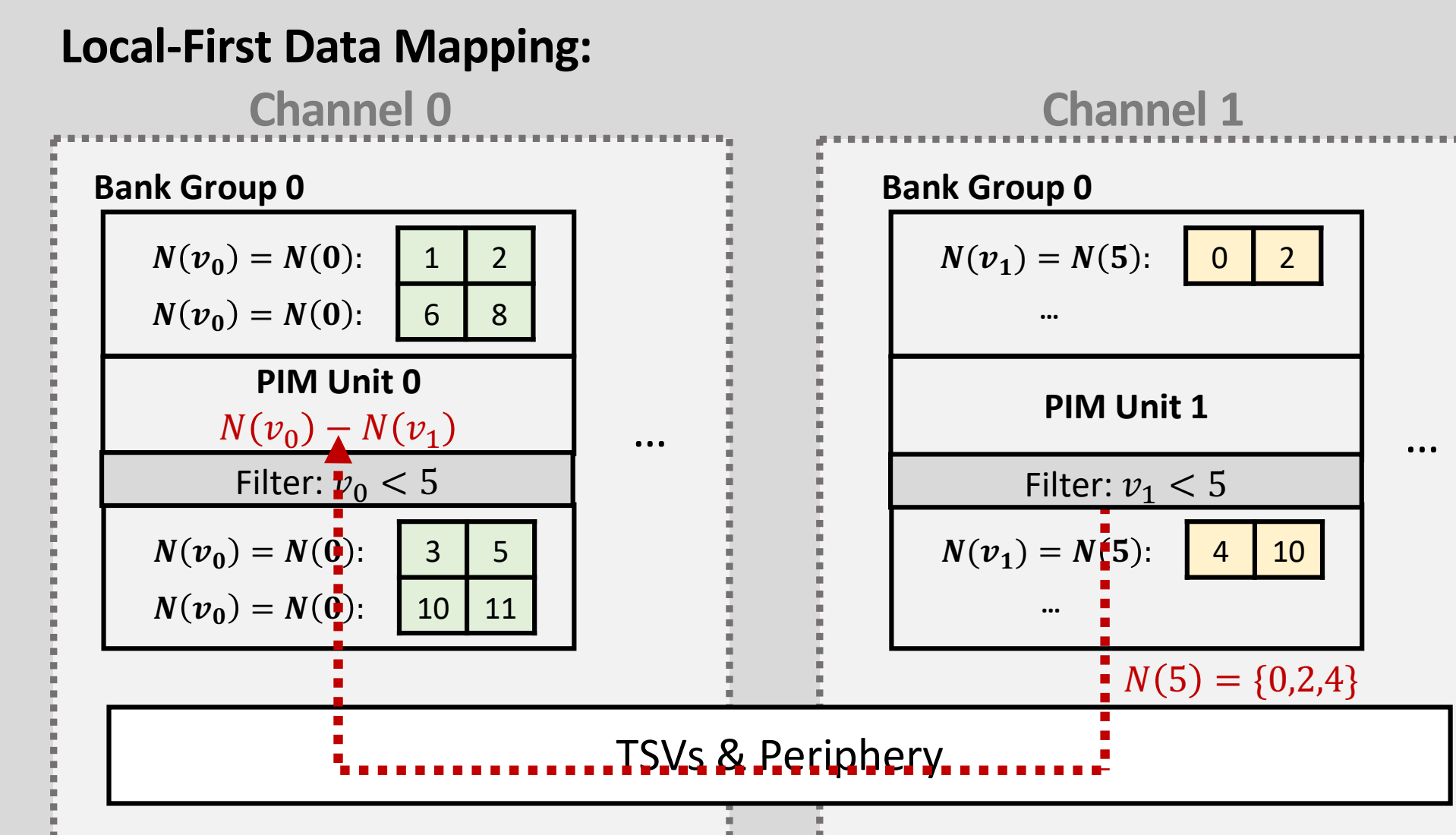
(b) HBM-DRAM Die and PIM unit memory accesses

## PIMMiner Framework Design

- **Conditional Access Filter (Filter)**: PIMMiner adds lightweight hardware dedicated to filter unnecessary data from memory. According to AutoMine, PIM units only need the nodes $v_2 < v_1$.

data $v_x$ from bank

threshold $th$

Restriction: $v_x < th$

subtractor 32b int

Filter Logic

condition $cmp: >, =, <$

$v_x - th$

send $v_x$ or $NULL$

- **Local-First Data Mapping (Remap):** We propose a new address mapping method to exploit low-latency and high-bandwidth local memory bank accesses for PIM units. The new mapping maps the data in the same neighbor list to the same bank group.

Local-First Data Mapping:

Channel 0
Bank Group 0
$N(v_0) = N(0)$: 1 2
$N(v_0) = N(0)$: 6 8
PIM Unit 0
$N(v_0) = N(v_1)$
Filter: $v_0 < 5$
$N(v_0) = N(0)$: 3 5
$N(v_0) = N(0)$: 10 11

Channel 1
Bank Group 0
$N(v_1) = N(5)$: 0 2
PIM Unit 1
Filter: $v_1 < 5$
$N(v_1) = N(5)$: 4 10
$N(5) = \{0,2,4\}$

TSVs & Periphery

- **Critical Data Duplication (Duplication):** To further reduce the remote memory accesses, PIMMiner stores critical data in the unused memory space in the memory banks.

- **Workload Scheduler with Stealing (Stealing):** PIMMiner achieves a PIM-aware workload scheduler by maintaining the execution tables and schedule tables on each PIM unit to address the load imbalance issues.

- **Working with specialized GPMI accelerator (Set):** PIMMiner can achieve even higher performance with integration of specialized GPMI accelerator in PIM units, such as the accelerator for intersection set operations.

## Evaluation Results

- Overall, by enabling all optimizations, PIMMiner achieves **15.91x** average speedup and **137.32x** maximum speedup over the baseline PIM.
- The performance improvement from PIMMiner optimizations:
  Filter: 2.01x average, 17.57x maximum speedup
  Remap: 1.39x average, 2.74x maximum speedup
  Duplication: 1.84x average, 3.05x maximum speedup
  Stealing: 3.01x average, 26.87x maximum speedup
  Set: 1.49x average, 2.26x maximum speedup

Base | Filter | Remap | Duplication | Stealing | Set

Normalized Time
5CC-CI 5CC-PP 5CC-AS 5CC-MI 5CC-YT 5CC-PA 5CC-LJ
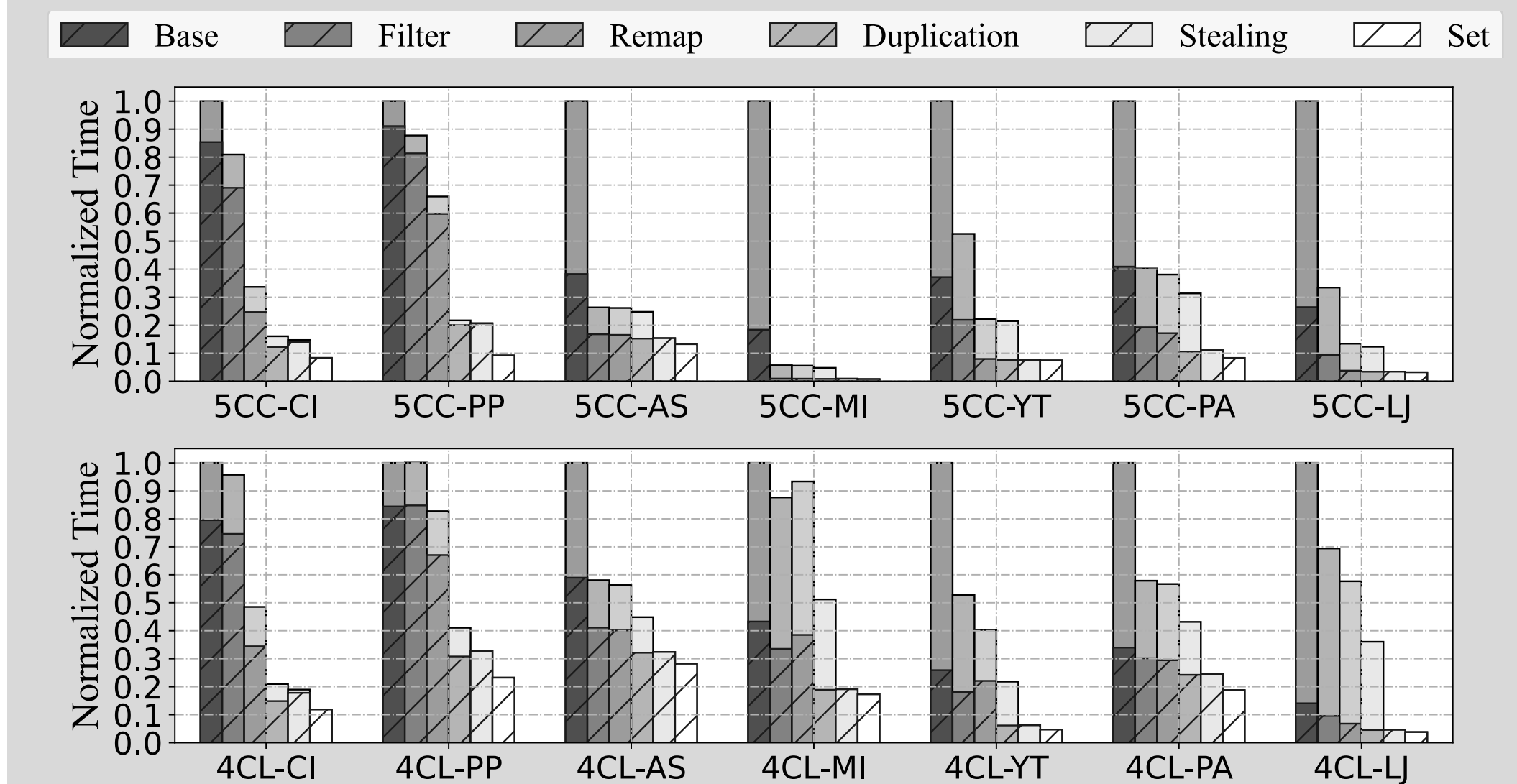4CL-CI 4CL-PP 4CL-AS 4CL-MI 4CL-YT 4CL-PA 4CL-LJ

Fig1: Performance of PIMMiner with the effectiveness of proposed optimizations. In each bar, we show the average time across cores (the solid line) and the total execution time(top of bar).

- Compare with other GPMI software systems:
  GraphPi [1]: **344x** average speedup
  AutoMine [2]: **109x** average speedup
- Compare with the hardware GPMI accelerators:
  Gramer [3]: **696x** average speedup
  FlexMiner [4]: **5.9x** average speedup
  DIMMining [5]: **1.3x** average speedup
  NDMiner [6]: **37x** average speedup

## References

[1] T. Shi et al., "Graphpi: High performance graph pattern matching through effective redundancy elimination," in SC, 20.
[2] D. Mawhirter et al, "Automine: harmonizing high-level abstraction and high performance for graph mining," in SOSP, 20.
[3] P. Yao et al., "A locality-aware energy-efficient accelerator for graph mining applications," in MICRO, 20.
[4] X. Chen et al., "Flexminer: a pattern-aware accelerator for graph pattern mining," in ISCA, 21.
[5] G. Dai et al., "Dimmining: pruning-efficient and parallel graph mining on near-memory-computing," in ISCA, 22.
[6] N. Talati et al., "Ndminer: accelerating graph pattern mining using near data processing." in ISCA, 22.

## Acknowledgement