PALP: Enabling and Exploiting Partition-Level Parallelism in Phase Change Memories

Shihao Song, Anup Das, Onur Mutlu and Nagarajan Kandasamy





Background: Bank Conflicts

PCM devices can serve multiple requests in parallel using bank-level parallelism



Requests to the same bank have to be served serially, known as bank conflict

Bank conflicts occur due temporal and spatial access

Key Observations

1st key observation: On average, 43% of PCM requests in SPEC 2017 workloads generate bank conflicts



Goals

- Analyze parallelism in PCM bank's partition through detailed circuit analysis
- Analyze the potential performance impact of enabling parallelism in a bank's partitions
- **Resolve read-read and read-write PCM bank** conflicts to improve PCM performance

Our Contributions

New mechanism to resolve read-write PCM bank

locality in a workload that lead to repeated access to multiple memory rows that map to the same bank

Memory bank conflicts reduce system performance by lowering bank utilization, causing CPU cores to stall

Decoupled Write Driver



2nd key observation: A PCM bank is implemented as a collection of partitions that operate mostly independently; sharing the sense amplifiers (to read) and the write drivers (to write)



The peripheral structures in PCM banks allow to read and program 128 PCM cells in parallel

conflicts using a new PCM command called **READ-WITH-WRITE (RWW)**

- Simple circuit modification to resolve read-read PCM
- bank conflicts using a new PCM command called **READ-WITH-READ (RWR)**
- New memory access scheduling to prioritize PCM requests that exploit a PCM bank's partition-level parallelism, over other requests, including the long outstanding ones

Enabling and Exploiting Partition-Level Parallelism in PCM

Performance Improvement vs. Baseline





| Operating | Write driver | circuit | Sense | Operations |
|-----------|--------------------|--------------|--------------|-------------------|
| Modes | Pulse shaper logic | Verify logic | Amplifier | Enabled |
| write | \checkmark | \checkmark | × | one write request |
| decoupled | × | \checkmark | \checkmark | two read requests |

serving a write and read request in our 2 PCM design (48 cycles)

Performance improvement 27%

serving two read requests in the new 2 PCM design (30 cycles)

Performance improvement 21%

PCM Commands

- **ACTIVATE(A):** activate the wordline and enable the access device for the PCM cells to be accessed
- **READ(R)/WRITE(W):** drive read or write current through the PCM cell. After this command executes, the data stored in the PCM cell is available at the output terminal of the sense amplifier, or the write data is programmed to the PCM cell
- **PRECHARGE(P)**: deactivate the wordline and bitline, and prepare the bank for the next access We introduce the following new commands
- **READ-WITH-WRITE (RWW):** connect the PCM bank's sense amplifiers and write drivers to the two decoded partitions
- **READ-WITH-READ (RWR):** connect the sense amplifiers and verify logic of the write drivers to the two decoded partitions
- **DECOUPLE (D)**: set M4 = OFF

| Evaluation | | | | | |
|---------------------------------------|---|--|--|--|--|
| Evaluation Methodology | MiBench workloads [20] | | | | |
| 1 Gem5 frontend to simulate ARMv8- | tiff2rgba, jpeg_decode, tiffdither, susan_smoothing, typeset | | | | |
| $\Lambda (aarch GA)$ with 0 corec | SPEC CPU2017 workloads [8] | | | | |
| A (ddicii04) With o Luies | cactusBSSN, bwaves, roms, parset, xz | | | | |
| 2. Hybrid DRAM-PCM memory system | Mixed (parallel) applications | | | | |
| 3. In-house cycle-level PCM simulator | AI-1 (4 copies each of deepsjeng, leela), AI-2 (4 copies each of mcf, exchange2), | | | | |
| for 8GB, 16GB, and 32GB PCM with | Visualization-1 (4 copies each of povray, blender), Visualization-2 (4 copies each of | | | | |
| DDR4 interface | povray, imagick), Scientific (4 copies each of cactusBSSN, bwaves) | | | | |

<u>Baseline</u>: Arjomand et al. "Boosting access parallelism to PCM-based main memory" in ISCA 2016 MultiPartition: Zhou et al., "An efficient parallel scheduling scheme on multi-partition PCM architecture" in ISLPED 2016





TRANSFER (T): sets M5 = ON and M6 = OFF

PCM Access Scheduling



RWW: Serving a read and a write request to RWR: Serving two read requests to different different partitions in the same bank partitions in the same bank

On average, execution time reduces by 51% vs. Baseline and 28% vs. MultiPartition

| Design | Technology Type | Vdd (V) | Critical Path Delay (ps) | Power (pJ/access |
|---------------|-----------------|---------|--------------------------|------------------|
| Baseline [37] | Double-gate | 0.9 | 1159.2 | 0.311 |
| Proposed PALP | Double-gate | 0.9 | 1453.2 | 0.364 |

On average, critical path delay increases by 25.3% On average, power consumption increases by 17%

> **Open source tool** https://github.com/drexel-DISCO/PALP

Both average and peak PCM power consumption of **PALP is within the RAPL limit** The average PCM power is at least 0.08pJ/access lower than the RAPL limit while the peak PCM power is at least 0.03pJ/access lower than the RAPL limit



On average, execution time reduces by 33% with DDR2 and 51% with DDR4