Exploiting Private Local Memories to Reduce the Opportunity Cost of Accelerator Integration

Emilio G. Cota Paolo Mantovani Luca P. Carloni

Columbia University



What accelerators, exactly?



Problem: Accelerators' Opportunity Cost

An accelerator is only of *utility* if it applies to the system's workload



Consequence: Integrating accelerators in general-purpose architectures is rarely cost-effective

Private Local Memories (PLMs)

Example: Sort Accelerator to sort FP vectors



Related Work

Key Observations:

1. Accelerators are mostly memory

4 An average of 69% of accelerator area is consumed by memory

Lyons et al., "The Accelerator Store", TACO'12

Accelerator examples: AES, JPEG encoder, FFT, USB, CAN, TFT controller, UMTS decoder..

2. Average Accelerator Memory Utilization is low Not all accelerators on a chip are likely to run at the same time

Related Work

Proposal [1]: The Accelerator Store
Shared memory pool that accelerators allocate from
Proposals [2,3]: Memory for Cache & Accelerators
Substrate to host cache blocks or accelerator buffers

Limitation: storage is external to accelerators

- Bigh-bandwidth PLMs cannot tolerate additional latency
- Somplicate accelerator designs
 - Hide PLM latency with pipelining, or reduce performance
- [1] Lyons et al., "The Accelerator Store: A Shared Memory Framework for Accelerator-Based Systems", TACO'12
- [2] Cong et al., "Bin: a Buffer-in-NUCA Scheme for Accelerator-rich CMPs", ISLPED'12
- [3] Fajardo et al., "Buffer-Integrated-Cache: a Cost-Effective SRAM Architecture for Handheld and Embedded Platforms", DAC'11

Our proposal: ROCA Observation #3: accelerator PLMs provide a de facto NUCA substrate

Goal:

To extend the LLC with PLMs when otherwise not in use



Applies to all accelerator PLMs, not only low-bandwidth ones

Minimal modifications to accelerators





ROCA. How to..

handle
 intermittent
 accelerator
 availability,



- 2 accommodate accelerators of different sizes,
- 3 transparently coalesce accelerator PLMs,

..with minimum overhead and complexity?

High-Level Operation

- core0's L1 misses on a read from 0xf00, mapped to the L2's *logical bank1*
- L2 bank1's tag array tracks block 0xf00 at acc2; sends request to acc2
- 3. acc2 returns the block to bank1
- 4. bank1 sends the block to core0



Solutional latency for hits to blocks stored in accelerators

- Return via the host bank guarantees the host bank is the only coherence synchronization point
 - No changes to coherence protocol needed

ROCA Host Bank



4-way example: 2 local, 2 remote ways

- Enlarged tag array for accelerator blocks
 - Ensures modifications to accelerators are simple
- Leverages Selective Cache Ways [*] to adapt to accelerators' intermittent availability
 - Dirty blocks are flushed to DRAM upon accelerator reclamation

Logical Bank Way Allocation



- Increasing associativity helps minimize waste due to uneven memory sizing across accelerators (Ex. 2 & 3)
- Power-of-two number of sets not required (Ex. 4), but
 - complicates set assignment logic [*]
 - requires full-length tags: modulo is not bit selection anymore

[*] André Seznec, "Bank-interleaved cache or memory indexing does not require Euclidean division", IWDDD'15

Coalescing PLMs



- PLM manager exports same-size dual-ported SRAM banks as multiported memories using MUXes
- ROCA requires an additional NoC-flit-wide port, e.g. 128b



- SRAMs are accessed in parallel to match the NoC flit bandwidth
 - Bank offsets can be computed cheaply with a LUT + simple logic
 - Discarding small banks and SRAM bits a useful option

ROCA: Area Overhead

- Host bank's enlarged tag array
 - 5-10% of the area of the data it tags (2b+tag per block)
- Tag storage for **standalone directory** if it wasn't there already
 - Inclusive LLC would require prohibitive numbers of recalls
 - Typical overhead: 2.5% of LLC area when LLC = 8x priv
- Additional logic: way selection, PLM coalescing logic
 - **Negligible** compared to tag-related storage

ROCA: Perf. & Efficiency

50



acc 128K	acc 192K	acc 128K	acc 192K	acc 192K	acc 192K	acc 128K
асс 192К	L2 256K	core	L2 256K	core	L2 256K	acc 192K
acc 192K	core	core	core	core	core	acc 128K
acc 192K	L2 256K	core	МС	core	L2 256K	acc 192K
acc 128K	core	core	core	core	core	acc 192K
acc 192K	L2 256K	core	L2 256K	core	L2 256K	acc 192K
acc 128K	acc 192K	acc 192K	acc 192K	acc 128K	acc 192K	acc 128K

Configurations:

- 2M S-NUCA baseline
- 8MB S-NUCA (not pictured)
- same-area 6M ROCA, assuming accelerators are 66% memory (below the typical 69%)

Assuming no accelerator activity,

- LLC MPKI reduction (%) 40 30 20 8M S-NUCA 6M ROCA 8M S-NUCA, gmean 6M ROCA, gmean 25 30 35 40 45 0 10 15 20 5 Workload 30 30 Performance improvement (%) 25 25 (%) 20 20 improvement 15 15 10 10 BIPJ -5 25 30 0 5 10 15 20 35 40 45 0 5 25 30 35 40 45 10 15 20 Workload Workload
- 6M ROCA can realize 70%/68% of the performance/energy efficiency benefits of a same-area 8M S-NUCA
 - retaining accelerators' potential orders-of-magnitude gains

Also in the paper

• Sensitivity studies sweeping accelerator activity over

- space (which accelerators are reclaimed)
- **time** (how frequently they are reclaimed)

- Key result: Accelerators with idle windows >10ms are prime candidates for ROCA
 - perf/eff. within 10/20% of that with 0% activity

core Lifetide Core Lifetide Core Core<

Accelerators can be highly-specialized, fixed-function units and still be of general-purpose utility

backup slides

Why Accelerators?

- Every generation provides less efficient transistors, i.e. power density is growing
- Single-threaded perf. improvements slowing down
- Parallelization gains bounded by Amdahl's Law

a.k.a. "the end of the multi-core era"

Esmaeilzadeh et al., "Dark Silicon and the End of Multicore Scaling", ISCA'11

Why Accelerators?

 If performance increases are to be sustained, we need efficiency gains well beyond what microarchitectural changes can provide

Accelerators achieve this via specialization

Based on adjusted SPECint® results 128 64 32 16 8 Intel Xeon 4 Intel Core Intel Pentium ▲ Intel Itanium Intel Celeron 2 AMD FX AMD Opteron ▲ AMD Phenom 1 AMD Athlon ▲IBM POWER • PowerPC $\frac{1}{2}$ ▲ Fujitsu SPARC Sun SPARC • DEC Alpha MIPS HP PA-RISC

Single-Threaded Integer Performance



Simulated Systems

cores	16 cores, i386 ISA, in-order IPC=1 except on memory accesses, 1GHz
L1 caches	Split I/D 32KB, 4-way set-associative, 1-cycle latency, LRU
L2 caches	8-cycle latency, LRU S-NUCA: 16ways, 8 banks ROCA: 12 ways
Coherence	MESI protocol, 64-byte blocks, standalone directory cache
DRAM	1 controller, 200-cycle latency, 3.5GB physical
NoC	5x5 or 7x7 mesh, 128b flits, 2-cycle router traversal, 1- cycle links, XY router
OS	Linux v2.6.34





Flushing delay

- 330 64-byte blocks (i.e. largest amount in our tests)
- sufficiently buffer DRAM controller
- 128b NoC flits

~10560 cycles, i.e. 10.5us at 1GHz