

PULP PLATFORM Open Source Hardware, the way it should be!

Working with RISC-V

Part 3 of 4 : PULP concepts

Luca Benini <lbenini@iis.ee.ethz.ch> Frank K. Gürkaynak <kgf@ee.ethz.ch>















Working with RISC-V

- Part 1 Introduction to RISC-V ISA
- Part 2 Advanced RISC-V Architectures

Part 3 – PULP concepts

- Cores
- Clusters
- Heterogeneous systems
- Accelerators



Part 4 – PULP based chips



Al Workloads from Cloud to Edge (Extreme?)







ETHZürich

4

RI5CY – Recap from Part 1

3-cycle ALU-OP, 4-cyle MEM-OP→IPC loss: LD-use, Branch







Hzürich

Nice – But what about the GOPS? Faster+Superscalar is not efficient! M7: 5.01 CoreMark/MHz-58.5 μW/MHz M4: 3.42 CoreMark/MHz-12.26 μW/MHz

5

ML & Parallel, Near-threshold: a Marriage Made in Heaven

- As VDD decreases. operating speed decreases
- However efficiency increases \rightarrow more work done per Joule
- Until leakage effects start to dominate
- (parMatrixMul2) [MOPs/mW] Put more units in parallel to get performance up and keep them busy with a Efficiency parallel workload

ML is massively parallel and scales well (P/S \uparrow with NN size)



Efficiency vs VDD chip01



Multiple RI5CY Cores (1-16)









7





ETH zürich



8

High speed single clock logarithmic interconnect



PEAC

ETHZürich

TCDM Interconnect

- Synchronous low latency crossbar based on binary trees
 - Single or dual channel (Each with 2ⁿ master ports, n=[0,1,2...]
 - Distributed Arbitration (Round robin)
 - Combinational handshake (single phase)
 - Deterministic access latency (1 cycle) + response in second cycle
 - Test&Set Supported

zürich

Word-level Interleaving

Emulates Multiported RAMs

- Slaves are pure memories.
- Bank conflicts can be alleviated \rightarrow higher Banking Factor (BF=1,2,4)



Peripheral Interconnect

Synchronous low latency crossbar based on binary trees

- Single or dual channel Each with 2n master ports, n=[0,1,2 ...]
- Distributed Arbitration (Round robin)
- Combinatorial handshake (single phase)
- Custom port mapping (Address ranges)
- Decoupled request and response path

Zürich

- Slaves have unpredictable latencies
- Used to build Peripheral systems
 - Slaves are pure generic peripherals like bridges, timers etc (Req/grant)
 - Mostly used to move data in and out from the cluster processors



FPU Interconnect (1/2)

- Synchronous low-power/low latency crossbar used to share several FPUs in a multicore system:
 - 2ⁿ master ports (n=[0,1,2 ...])
 - 2^m slave ports (m=[0,1,2 ...]) → General purpose FPUs(ADD/SUB, MUL etc)
 - Combinatorial handshake (single phase)
 - Allocator
 - Random: given a request, a random fpu is choosen
 - Optimal: Maxime the utilization of FPUs





FPU Interconnect

- Features:
 - Independent response paths for each sub-FPU block
 - fully pipelined FPU sub-blocks with different latencies
 - Two Operators (A,B), one command (OP) and the ID are carried to the FPU.
 - No need of Flow control on the FPU side.
 - Flexible and parametrizable



Example of FPU attached to the FPU interconnect







DMA for data transfers from/to L2



EHzürich

14

PULP MCHAN DMA Engine

- A DMA engine optimized for integration in tightly coupled processor clusters
 - Dedicated, per-core non blocking programming channels
 - Ultra low latency programming (~10 cycles)
 - Small footprint (30Kgates) : avoid usage of large local FIFOs by forwarding data directly to TCDM (no store and forward)
 - Support for multiple outstanding transactions
 - Parallel RX/TX channels allow achieving full bandwidth for concurrent load and store operations

Configurable parameters:

of core channels

H zürich

- Size of command queues
- Size of RX/TX buffer
- # of outstanding transactions







Fast synchronization and Atomics



ETHZürich

16

Synchronization & Events



PULP Cluster Event Unit



ETHZürich

18

Energy-efficient Event Handling





19





- Fully parallel access to SCU: Barrier cost constant
- Primitive energy cost: Down by up to 30x
- Minimum parallel section for 10% overhead in terms of ...
 - ... cycles: ~100 instead of > 1000 cycles
 - ... energy: \sim 70 instead of > 2000 cycles

Results: Mutex



- Sequential execution: Cycle overhead always large
- TAS-variable inherently well-suited for mutex; lower cycle savings compared to barrier

- SCU still avoids L1 accesses: Energy of TAS mutex up to 1.6x higher
- Smallest parallel section for 10% energy overhead:
 - ~1000 instead of 1600 cycles



How do we work: Initiate a DMA transfer



ETHZürich

THE TOTAL

25.7.2018

22

Data copied from L2 into TCDM



ETH zürich

TWITE DES

25.7.2018

23

Once data is transferred, event unit notifies cores



ETHZürich



25.7.2018

24

Cores can work on the data transferred



ETHZürich



25.7.2018

25

Once our work is done, DMA copies data back



ETHZürich

25.7.2018

26

During normal operation all of these occur concurrently



ETHZürich

25.7.2018

Shared instruction cache with private "loop buffer"



ETHZürich

28

ULP (NT) Bottleneck: Memory

- "Standard" 6T SRAMs:
 - High VDDMIN
 - Bottleneck for energy efficiency
 - >50% of energy can go here!!!
- Near-Threshold SRAMs (8T)
 - Lower VDDMIN
 - Area/timing overhead (25%-50%)
 - High active energy
 - Low technology portability

Standard Cell Memories:

- Wide supply voltage range
- Lower read/write energy (2x 4x)
- High technology portability
- Major area overhead $4x \rightarrow 2.7x$ with controlled placement



ETHZürich

I\$: a Look Into 'Real Life' Applications

SCM-BASED I\$ IMPROVES EFFICIENCY BY ~2X ON SMALL BENCHMARKS, BUT... Applications on PULP

Survey of State of The Art 1,2 Latch based I\$ **Exixting ULP processors** Degradation REISC (ESSCIRC2011) 64b Sleepwalker (ISSCC 2012) 128b **9**,0,6 Bellevue (ISCAS 2014) 128b **E** 0,4 SHORT JUMP LOOP **BASED APPLICATIONS** BFS ----*--- MD **0**,2 CT ----- FAST LONG JUMP APPLICATIONS SLIC -HOG LIBRARY BASED Issues: 522 250 1 at Area Overhead of SCMs (4Kb/core not affordable....) 1) I\$Size [B]

2) Capacity miss (with small caches)

ETH zürich

3) Jumps due to runtime (e.g. OpenMP, OpenCL) and other function calls



- Shared instruction cache
 - OK for data parallel execution model
 - Not OK for task parallel execution model, or very divergent parallel threads

Architectures

- SP: single-port banks connected through a read-only interconnect
 - Pros: Low area overhead
 - Cons: Timing pressure, contention
- MP: Multi-ported banks
 - Pros: High efficiency
 - Cons: Area overhead (several ports)

Results

- Up to 40% better performance than private I\$
- Up to 30% better energy efficiency
- Up to 20% better energy*area efficiency



ETHZürich

31

Results: RV32IMCXpulp vs RV32IMC

- 8-bit convolution
 - Open source DNN library
- 10x through xPULP
 - Extensions bring real speedup
- Near-linear speedup
 - Scales well for regular workloads.
- 75x overall gain



An additional I/O controller is used for IO



ETH zürich











- Efficient use of system resources
- HW support for double buffering allows continuous data transfers
- Multiple data streams can be time multiplexed





PULP interrupts controller (INTC)

- It generates interrupt requests from 0 to 31
- Mapped to the APB bus
- Receives events in a FIFO from the SoC Event Generator (i.e. from peripherals)
 - Unique interrupt ID (26) but different event ID
- Mask, pending interrupts, acknowledged interrupts, event id registers
- Set, Clear, Read and Write operations by means of load and store instructions (memory mapped operations)
- Interrupts come from:
 - Timers

ETH zürich

- GPIO (rise, fall events)
- HWCE
- Events i.e. uDMA



Tightly-coupled HW Compute Engine



42

Hardware Processing Engines (HWPEs)



PEAC

43

ETH zürich





externally, uses memory accesses (master ports)



ETH zürich



Peripheral access (target port)



25.7.2018

HW Convolution Engine



F. Conti and L. Benini, "A ultra-low-energy convolution engine for fast brain-inspired vision in multicore clusters," *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2015, pp. 683-688.

ETHZürich



45

HWCE Sum-of-Products



Need µW-range always-on Intelligence



47

HD-Based smart Wake-Up Module





PEAC

48

ETH zürich



HD-Based smart Wake-Up Module



ETHZürich



HD-Based smart Wake-Up Module





Results (post-P&R) - TO done	
Technology	GF22 UHT
Area	670kGE
Max. Frequency	3 MHz
SCM-Memory	32 kBit
Module Power Consumption (@ 1 kSPS/channel, 3 channels)	~ 15uW



PULP includes Cores+Interco+IO+HWCE → Open Platform





52

P U P

Nice, but what exactly is "open" in Open Source HW?

- Only the first stage of the silicon production pipeline can be open HW
 → RTL source code (in an HDL such as SystemVerilog)
- Later stages contain closed IP of various actors + tool licensing issues



0 SIC **B**







Parallel Ultra Low Power

Luca Benini, Davide Rossi, Andrea Borghesi, Michele Magno, Simone Benatti, Francesco Conti, Francesco Beneventi, Daniele Palossi, Giuseppe Tagliavini, Antonio Pullini, Germain Haugou, Manuele Rusci, Florian Glaser, Fabio Montagna, Bjoern Forsberg, Pasquale Davide Schiavone, Alfio Di Mauro, Victor Javier Kartsch Morinigo, Tommaso Polonelli, Fabian Schuiki, Stefan Mach, Andreas Kurth, Florian Zaruba, Manuel Eggimann, Philipp Mayer, Marco Guermandi, Xiaying Wang, Michael Hersche, Robert Balas, Antonio Mastrandrea, Matheus Cavalcante, Angelo Garofalo, Alessio Burrello, Gianna Paulin, Georg Rutishauser, Andrea Cossettini, Luca Bertaccini, Maxim Mattheeuws, Samuel Riedel, Sergei Vostrikov, Vlad Niculescu, Hanna Mueller, Matteo Perotti, Nils Wistoff, Luca Bertaccini, Thorir Ingulfsson, Thomas Benz, Paul Scheffler, Alessio Burello, Moritz Scherer, Matteo Spallanzani, Andrea Bartolini, Frank K. Gurkaynak,

and many more that we forgot to mention

http://pulp-platform.org



@pulp_platform