

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

Open Platforms for Energy-Efficient Scalable Computing

Luca Benini < Ibenini@iis.ee.ethz.ch,Iuca.Benini@unibo.it>



Computing is Power Bound: HPC





2

Computing is Power Bound: 5G/6G



Baseband station processing: 10x every 3 years!



Computing is Power Bound: ML

Largest datacenter <150MW



Machine Learning (training): 10x every 2 years



Technology Scaling?

ETH zürich

TSMC, ISSCC21



Energy Efficiency
$$\left(\frac{1}{Power \cdot Time}\right) \rightarrow 10x \text{ every } 12 \text{ years...}$$

Efficient Architecture: Heterogeneous+Parallel



Heterogeneous + Parallel... Why?

Processors can do two kinds of useful work:

Decide (jump to different program part)

- Modulate flow of instructions
- Mostly sequential decisions:
 - Don't work too much
 - Be clever about the battles you pick (latency is king)
- Lots of decisions
 Little number crunching

Compute (plough through numbers)

- Modulate flow of data
- Embarassing data parallel:
 - Don't think too much
 - Plough through the data (throughput is king)
- Few decisions
 Lots of number crunching
- Today's workloads are dominated by "Compute":
 - Tons of data, few (as fast as possible) decisions based on the computed values,
 - "Data-Oblivious Algorithms" (ML, or better DNNs are so!)
 - Large data footprint + sparsity

ETH zürich

How to design an efficient "Compute" fabric?

Compute Efficiency: D (...and I) Movement is Key





PE: Snitch, a Tiny RISC-V Core

A versatile building block



Simplest core: around 20KGE

- Speed via simplicity (1GHZ+)
- L0 Icache/buffer for low energy fetch
- Shared L1 for instruction reuse (SPMD)

■ Extensible → "Accelerator" port

- Minimal baseline ISA (RISC-V)
- Extensibility: Performance through ISA extensions (via accelerator port)

■ Latency-tolerant → Scoreboard

- Tracks instruction dependencies
- Much simpler than OOO support!



F. Zaruba, F. Schuiki, T. Hoefler and L. Benini, "Snitch: A Tiny Pseudo Dual-Issue Processor for Area and Energy Efficient Execution of Floating-Point Intensive Workloads," in *IEEE Transactions on Computers*, vol. 70, no. 11, pp. 1845-1860, 1 Nov. 2021

Snitch PE: ISA Extension for efficient "Compute"

- How can we remove the Von Neumann Bottleneck?
- Targeting "compute" code



Memory access, operation, iteration control – can we do better? Note: memory access (>1 cycle even for L1) \rightarrow need latency tolerance for LD/ST



Stream Semantic Registers

LD/ST elision

- Intuition: High FPU utilization ≈ high energy-efficiency
- Idea: Turn register read/writes into implicit memory loads/stores.
- Extension around the core's register file
- Address generation hardware

scfg 0, %[a], ldA loop: scfg 1, %[b], ldB loop: fmadd r2, ssr0, ssr1

- Increase FPU/ALU utilization by ~3x up to 100%
- SSRs ≠ memory operands
 - Perfect prefetching, latency-tolerant
 - 1-3 SSR (2-3KG/SSR)





Floating-point Repetition Buffer

Remove control flow overhead in compute stream



- Programmable micro-loop buffer
- Sequencer steps through the buffer, independently of the FPU
- Integer core free to operate in parallel:
 Pseudo-dual issue
- High area- and energy-efficiency

```
      mv
      r0, zero
      frep
      r1, 1

      loop:
      ______
      loop:
      ______
      loop:

      addi
      r0, 1
      ______
      loop:
      ______

      fmadd
      r2, ssr0, ssr1
      fmadd
      r2, ssr0, ssr1

      bne
      r0, r1, loop
      ______
      fmadd
      r2, ssr0, ssr1
```



RISC-V ISA Extension for Target Workload

Mixed precision

Efficient DNN inference & training



Inference ≠ Training Quantization

- Inference: INT8 quantization is SoA
- Training: High dynamic range needed for weights and weight updates

fp32 is still standard for DNN training workloads. Low precision training with **bf18** and **fp8**

Support a wide variety of FP formats and instructions:

- Standard: **fp64**, **fp32**, **fp16**, **bf16**
- Low precision: fp8, altfp8
 - fp8 (1-4-3): forward prop.
 - altfp8 (1-5-2): backward prop.
 - Exp. ops: accumulation



FPU and SDOTP unit





EXFMA vs EXSDOTP

Expanding Dot Product w. accumulation

FPU Vectorial FMA							FPU SDOTP					
[FP16	FP16	FP16	FP16	rs0		FP16	FP16	FP16	FP16	rs0	
[FP16	FP16	FP16	FP16	rs1		FP16	FP16	FP16	FP16	rs1	
FP32		FP32		rs2		FP32		FP32		rs2		
	FP32		FP32		rd		FP32		FP32		rd	

- $rd_{32} = rs\theta_{16} * rs1_{16} + rs2_{32}$
- Expanding FMA: Unbalanced
- Consumes only half of rs0, rs1 and produces the whole rd

- $rd_{32} = rs0_{0,16} * rs1_{0,16} + rs0_{1,16} * rs1_{1,16} + rs2_{32}$
- Expanding SDOTP: Balanced
- The whole rs0, rs1, rs2, rd are used



Cascade of EXFMAs vs EXSDOTP



Non-distributive FP addition \rightarrow **Precision Loss**

- Fused EXSDOTP (i.e. lossless)
- Single normalization and rounding step
- Smaller area and shorter critical path
- Product by-pass to compute fused three-term addition (vector inner sum)
- Stochastic rounding supported (+3% area)



What About Sparsity? Sparse SSR Streamer



Sparse SSR Functional Benefits

- Enable continuous useful issues in sparse LA
- Accelerate uses beyond sparse LA
 - Dense workloads: dense LA, CONV, FFT, ...
 - ISSRs as dynamic pattern streamers: stencil codes, codebook decoding, ...
 - Index matching in graph operations

Provide SoA capabilities without common compromises

- Indirection, intersection, and union
- Efficient inflight *or* AoT data transform
- No sparsity or structural requirements (e.g. A100¹)
- No imposed dataflow (e.g. OuterSpace²)

```
dotp:call
             conf ft0 mst match
   call
            conf_ft1_mst_match
   frep.s
   fmadd.d fa1, ft0, ft1, fa1
vadd:call
             conf_ft0_mst_mergecompute
                                 only
           conf_ft1_mst_merge
   call
   call
           conf_ft2_slv_wb
   frep.s
           ft2, ft0, ft1
   fadd.d
```



"2:4 structured pattern"



Sparse SSR Performance Benefits

- Notable single-core speedups over RV baseline
 - CsrMV: up to 7.0× faster, 79% FP util.
 - *SpV*+*SpV*: up to **9.8**× faster / higher FP util.
 - *SpV*·*SpV*: up to **7.7**× faster / higher FP util.
 - VTI (3D stencil code): up to 2.9× faster, 78% FP util.

Significant benefits in multicore cluster:

- CsrMV : up to 5.0× faster, 2.9x less energy
- CsrMSpV : up to 5.8× faster, 3.0x less energy
- VTI: up to 2.7× faster
- Notably higher peak FP utilizations than SoA CPUs (69×), GPUs (2.8×) on CsrMV





Efficient PE (snitch) architecture in perspective

- **1.** Minimize control overhead \rightarrow Simple, shallow pipelines
- **2.** Reduce VNB \rightarrow amortize IF: SSR-FREP + SIMD (Vector processing)
- 3. Hide memory latency \rightarrow non-blocking (indexed) LD/ST+dependency tracking
- 4. Highly expressive, domain-specific instruction extensions (thanks, RISC-V!)





Compute Efficiency: the Cluster (PEs + On-chip TCDM)





The Cluster: Design Challenges

Efficient PE

- Hide TCDM "residual" latency
- Remove Von Neumann Bottleneck

Low latency access TCDM

- Multi-banked architecture
- Fast logarithmic interconnect

Fast synchronization

- Atomics
- Barriers





High speed logarithmic interconnect



A. Rahimi, I. Loi, M. R. Kakoee and L. Benini, "A fully-synthesizable single-cycle interconnection network for Shared-L1 processor clusters," 2011 Design, Automation & Test in Europe, 2011, pp. 1-6.23

Efficient Explicit Global Data Mover

hide L2main memory latency



- 64-bit AXI DMA explicit double-buffered transfers – better than D\$
- Tightly coupled with Snitch (<10 cycles configuration)
- Operates on wide 512-bit data-bus
- Hardware support to copy 2-4-dim shapes
- Higher-dimensionality handled by SW
- Intrinsics/library for easy programming
- Domain-specific autotilers
 ETH zürich

. . .

```
// setup and start a 1D transfer, return transfer ID
uint32_t __builtin_sdma_start_oned(
    uint64_t src, uint64_t dst, uint32_t size, uint32_t cfg);
// setup and start a 2D transfer, return transfer ID
uint32_t __builtin_sdma_start_twod(
    uint64_t src, uint64_t dst, uint32_t size,
    uint32_t sstrd, uint32_t dstrd, uint32_t nreps, uint32_t cfg);
// return status of transfer ID tid
uint32_t __builtin_sdma_stat(uint32_t tid);
// wait for DMA to be idle (no transfers ongoing)
void __builtin_sdma_wait_for_idle(void);
```

Snitch Cluster Architecture





Where does the Energy go?



Baikonur

The first Snitch-based test chip.



> 50% Snitch FPU Seq Regs Rest FPSS 188 142 21 13 8 4 Regfile SSR LSU Rest Regfile Perf Cnt LSU Care FPU 142 21 24 16 2 5 21 11 96 8 kGE

Snitch compute cluster:

- 8x RV32G Snitch cores
- 8x large FPUs
- > 40% energy spent on FPU
- FPU utilization >80%+ (10% banking conflicts)
- High energy-efficiency of ~80 Gdpflop/s/W
- 104 Gspflop/s/W similar to accelerators
- 9 mm² prototype in GF 22FDX
- Testbed for key architectural components
 - Snitch octa-core cluster
 - Ariane cores

Very efficient, versatile, compute cluster! How do we scale?

Efficient Cluster architecture in perspective

- **1.** Memory pool efficient sharing of L1 memory
- **2.** Fast and parsimonious synchronization
- **3.** Data Mover + Double buffering explicitly managed block transfers at the boundary
- 4. More cores and more memory per cluster... that would be nice!



Compute Efficiency: the Chip(let) (Clusters+Off-die Mem)





From Snitch Cluster

8 Snitch compute cores

Single-stage, small Integer control core

9th Core: DMA

- 512 bit data interface
- Efficient data movement

128 kB TCDM

- Scratchpad for predictable memory accesses
- 32 Banks

Custom ISA extensions

- Xfrep, Xssr
- New: Xissr sparsity support

1 FPU per Snitch core

- Decoupled and heavily pipelined
- Multi-format FPU (+SIMD)
- New: Minifloat support + SDOTP



to Group AXI Wide 512bit





To Snitch Group

to Group AXI Narrow 64bit



to Group AXI Wide 512bit









From Snitch Group



to Toplevel AXI Narrow 64bit

to Toplevel AXI Wide 512bit



-

To Multi-Group



Occamy – Chiplet Architecture





Occamy NoC: Efficient and Flexible Data Movement



[1] Burrello, Alessio, et al. "Dory: Automatic end-to-end deployment of real-world dnns on low-cost iot mcus." IEEE Transactions on Computers 70.8 (2021): 1253-1268.

Problem: HBM Accesses are critical in terms of

- Access energy
- Congestion
- High latency

Instead reuse data on lower levels of the memory hierarchy

- Between clusters
- Across groups

Smartly distribute workload

- Clusters: DORY framework for tiling strategy [1]
- Chiplets: E.g. Layer pipelining

Big trend!





Occamy's C2C Link







C2C Link

Scalable and fault tolerant Chiplet2Chiplet Link



Network Layer

- Full AXI4 interface
- AXI4 to AXI stream converter

Data Link Layer

- Credit-based flow control
- RX synchronisation

Channel Allocator

- Chops and reshuffles payload
- Fault tolerance mechanisms

Physical layer (decoupled)

- Source-synchronous DDR sampling
- Scalable number of channels (38 x 8-bit full-duplex DDR channels in Occamy)
- Compatible with fast serial phy (e.g. HBI)



Occamy Chiplet



- GF12, target 1GHz (typ)
- 2 AXI NoCs (multi-hierarchy)
 - 64-bit
 - 512-bit with "interleaved" mode
- Peripherals
- Linux-capable manager core CVA6
- 6 Quadrants: 216 cores/chiplet
 - 4 cluster / quadrant:
 - 8 compute +1 DMA core / cluster
 - 1 multi-format FPU / core (FP64,x2 32, x4 16/alt, x8 8/alt)
- 8-channel HBM2e (8GB) 512GB/s
- D2D link (Wide, Narrow) 70+2GB/s
- System-level DMA
- SPM (2MB wide, 512KB narrow)

ETHzürich Peak 384 GDPflop/s per chiplet – Taped-out in July (fully OSHW)

Efficient Chiplet architecture in Perspective

- **1.** Multi-cluster single-die scaling \rightarrow strong latency tolerance, modularity
- **2.** NoC for flexible Clus2Clus, Clus2Mem traffic \rightarrow reduce pressure to Main memory
- **3.** Top level NoC Routes to "local main memory" / "global main memory" balanced BW
- 4. Modular chiplet architecture: HBM2e, NoC-wrapped C2C, multi-chiplet ready



Back to the cluster... Can we make it Bigger?

Why?

- Better global latency tolerance if L1 > 2*Latency*Bandwidth
- Easier to program (data-parallel, functional pipeline...)
- Smaller data partitioning overhead

An efficient many-core cluster with low-latency shared L1

- 256+ cores
- 1+ MiB of shared L1 data memory
- ≤ 5 cycles latency (without contention)

Physical-aware design

- WC Frequency > 500 Mhz
- Targeting iso-frequency with small cluster





Hierarchical Physical Architecture

Tile

- 4 32-bit cores
- 16 banks
- Single cycle memory access

Group

- 64 cores
- 256 banks
- 3-cycles latency

Cluster

- 256 cores
- 1 MiB of memory (1024 banks)
- 5-cycles latency



ETHzürich TopH: Butterfly Multi-stage Interconnect 0.3req/core/cycle

Benchmarks: can we compete with the impossible?

- Baseline: idealized Top_X
 - Fully connected logarithmic crossbar between 256 cores and 1024 banks
- Cycle-accurate RTL simulation
 - Matmul
 - Multiplication of two 64 × 64 matrices
 - 2dconv
 - 2D Convolution with a 3 × 3 kernel
 - dct
 - 2D Discrete Cosine Transform on 8 × 8 blocks in local memory
- Top_H has a performance penalty of at most 20%, on all kernels





Energy Analysis (500 MHz, TT, 0.80 V, 25 °C)

Breakdown of the energy consumption per instruction:





- About half of it, 4.5pJ, by the interconnect
- Remote loads: twice the energy of a local load
 - Despite crossing the whole cluster, twice!

ETHzürich **Terapool (1024 cores) seems viable** @ 7 cycles latency



MemPool-3D

- Memory-on-Logic implementation of MemPool
- Implementation of MemPool with 1, 2, 4, 8 MiB of L1
- Leading to a higher utilization of the memory die





Groups: MemPool-2D vs. MemPool-3D



Similar trends for Wire length and #Buffers



Conclusion

- Energy efficiency quest: PE, Cluster, SoC
- Key ideas
 - Deep PE optimization → extensible ISAs (RISC-V!)
 - VNB removal + Latency hiding: large OOO processors not needed
 - Low-overhead work distribution → large "mempool"
 - Heterogeneous architecture →host+accelerator(s)
- Game-changing technologies
 - "Commoditized" chiplets: 2.5D, 3D
 - Computing "at" memory (DRAM mempool)
 - Coming: optical IO and smart NICs, swiches
- Challenges:
 - RV Host?
 - RV HPC software ecosystem?





High Bandwidth / High Memory Capacity General-Purpose Many-Core CPU High Bandwidth SRAM + Large Capacity DRAM or NVM



Monte Cimone: the first RISC-V Cluster







Designed for HPC "pipe cleaning"

RISC-V Challenge... SW ecosystem (?)

In 2021 we designed and built **Monte Cimone**, the **first physical prototype** and test-bed of a **complete RISC-V (RV64) compute cluster** integrating **compute, interconnect**, *a complete software stack for HPC* and a full-featured **system monitoring** *infrastructure*.

- 1. Ported and assessed the maturity of a HPC software stack composed of:
 - SLURM job scheduler, NAS filesystem, Spack package manager
 - compilers toolchains, scientific and communication libraries,
 - a set of HPC benchmarks and applications,
 - ExaMon datacenter automation and monitoring framework.



- Characterized the HPL and STREAM benchmarks w. the toolchain and libraries installed by Spack: Compared against Armida (ARMv8a ThunderX2) and Marconi100 (ppc64le, IBM Power9)
- 3. Measured and analyzed MC power and energy using Examon



Monte Cimone Hardware



4x E4 RV007 1U Custom Server Blades:

- 2x SiFive U740 SoC with 4x U74 RV64GCB cores
- 16GB of **DDR4**
- 1TB node-local NVME storage
- PCle expansion card w/InfiniBand HCAs
- Ethernet + IB parallel networks

SiFive U740 SoC w. 7 separated power rails:

- Core complex, IOs, PLLs, DDR subsystem and PCIe one.
- Board implements distinct shunt resistors



Performance Characterization

Monte Cimone vs ARMv8a, PPC64le:

HPL and Stream benchmarks on two SoA computing nodes:

- <u>Marconi100 (</u>ppc64le, IBM Power9)
- <u>Armida (</u>ARMv8a, Marvell ThunderX2)
- Same benchmarking boundary conditions

ETH zürich

- Vanilla unoptimized
 libraries
- software stack deployed via SPACK package manager



Monte Cimone vs ARMv8a, ppc64le

Power Characterization



	Idle		HPL		STREAM.L2		STREAM.DDR		QE		Boot	
Line											R1	R2
	[mW]	[%]	[mW]	[%]	[mW]	[%]	[mW]	[%]	[mW]	[%]	[mW]	[mW]
COLE	3075	64	4097	69	3714	68	3287	62	3825	67	984	2561
ddr_soc	139	3	177	3	170	3	232	4	176	3	59	197
io	20	0	20	0	20	Q	20	0	20	0	5	20
pll	1	0	1	0	1	0	1	0	1	0	0	2
pcievp	521	11	527	9	524	10	522	10	530	9	12	231
pcievph	555	12	554	9	554	10	555	10	561	10	1	395
ddr_mem	404	8	440	7	401	7	592	11	434	8	275	467
ddr_pll	28	1	28	1	28	1	28	1	28	1	0	29
ddr_vpp	67	1	90	2	73	1	98	2	95	2	49	122
Total	4810	100	5935	100	5486	100	5336	100	5670	100	1385	4024

Core complex @ boot process:

- 0.981W of leakage only power (32% of the Idle power)
- 0.514W O. S. idle power (17% of the Idle power)
- 1.577W of dynamic and clock tree power (51% of the idle power).



Luca Benini, Alessandro Capotondi, Alessandro Ottaviano, Alessandro Nadalini, Alessio Burrello, Alfio Di Mauro, Andrea Borghesi, Andrea Cossettini, Andreas Kurth, Angelo Garofalo, Antonio Pullini, Arpan Prasad, Bjoern Forsberg, Corrado Bonfanti, Cristian Cioflan, Daniele Palossi, Davide Rossi, Davide Nadalini, Fabio Montagna, Florian Glaser, Florian Zaruba, Francesco Conti, Frank K. Gürkaynak, Georg Rutishauser, Germain Haugou, Gianna Paulin, Gianmarco Ottavi, Giuseppe Tagliavini, Hanna Müller, Lorenzo Lamberti, Luca Bertaccini, Luca Valente, Luca Colagrande, Luka Macan, Manuel Eggimann, Manuele Rusci, Marco Guermandi, Marcello Zanghieri, Matheus Cavalcante, Matteo Perotti, Matteo Spallanzani, Mattia Sinigaglia, Michael Rogenmoser, Moritz Scherer, Moritz Schneider, Nazareno Bruschi, Nils Wistoff, Pasquale Davide Schiavone, Paul Scheffler, Philipp Mayer, Robert Balas, Samuel Riedel, Sergio Mazzola, Sergei Vostrikov, Simone Benatti, Stefan Mach, Thomas Benz, Thorir Ingolfsson, Tim Fischer, Victor Javier Kartsch Morinigo, Vlad Niculescu, Xiaying Wang, Yichao Zhang, Yvan Tortorella, all our past collaborators and many more that we forgot to mention

http://pulp-platform.org



@pulp_platform

Fast synchronization and Atomics

Atomic instructions

- Leverage fast memory access
- Based on per-bank atomic adapters

Fast synchronization

- Atomics for generic primitives
- Accelerating barriers in HW (make the common case fast)
- Minimize idle power while waiting





Hardware-Accelerated, Event-Based Barrier



Barrier: Results



- 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: ~70 instead of > 2000 cycles



Stencils: Performance Results

- Kernel 1 (16x4x2 output tile, radius 3, 12 timesteps)
 - Single-core FPU utilization: **77%**
 - 8-core parallel speedup: 6.5× (ideal limit 8×)
- Kernel 2 (3x3x3 output tile, radius 3, 12 timesteps)
 - Single-core FPU utilization: 83%
 - 8-core parallel speedup: **5.4**× (ideal limit 6.75×)









Add date or a third information here

58

High Single-Core Efficiency

- Instructions almost entirely useful FP
 - Pseudo dual-issue hides SSR config
- Remaining stalls:
 - 1st point: cold I\$, regs
 - RAW dependencies (out of registers)
 - In-order stores (out of SSRs)
 - ISSR index bottleneck (limited buffers)
 - PDI knock-on (not a real stall)

ETH zürich

761 add s0, s9, a5	822 fsub.d fa5, ft0, ft1	874 fmadd.d ft7, fs10, fa1, ft7	922 fmadd.d fs1, fa6, fs7, fs1	972 fadd.d ft3, ft3, ft5
762 addi a4, s0, 912	823 fmadd.d fs0, fa6, fs6, fs0	875 fmadd.d ft8, fs11, fa1, ft8	923 fmadd.d fs2, fa6, fs8, fs2	973 fadd.d ft4, ft4, ft8
763 scfgwi	824 fmadd.d fs1, fa6, fs7, fs1	876 fld fa1, 8(t0)	924 fmadd.d fs3, fa6, fs9, fs3	974 fmadd.d ft5, fa2, ft0, fa1
764 lw s1, 0(s3)	825 fmadd.d fs2, fa6, fs8, fs2	877 fld fa2, 16(t0)	925 fmadd.d fs4, fa6, fs10, fs4	975 fadd d fs8 fs8 fs0
•••	<pre>826 fmadd.d fs3, fa6, fs9, fs3</pre>	878 fsub.d fs6, ft0, ft1	926 fmadd.d fs5, fa6, fs11, fs5	0.76 framid ft ft ft ft
768 scfgwi	827 fmadd.d fs4, fa6, fs10, fs4	879 fsub.d fs7, ft0, ft1	927 fnmsub.d fs0, fa6, ft9, fs0	970 (Sgill, 0 (10), (11), (11)
769 lw a4, 0(s8)	828 fmadd.d fs5, fa6, fs11, fs5	880 fsub.d fs8, ft0, ft1	928 fmadd.d fs1, fa6, ft10, fs1	977 TSUD.U TLS, TLS, TLS
	829 fnmsub.d fs0, fa6, ft9, fs0	881 fsub.d fs9, ft0, ft1	929 fmadd.d fs2, fa6, ft11, fs2	978 fmadd.d ft7, fa2, ft0, fa1
7/3 sctgw1	830 +madd.d +s1, +a6, +t10, +s1	882 fsub.d fs10, ft0, ft1	930 tmadd.d ts3, ta6, ta3, ts3	979 fmsub.d ft8, fa2, ft1, ft0
//6 +10 +a0, 8(sp)	831 fmadd.d fs2, fa6, ft11, fs2	883 fsub.d fs11, ft0, ft1	931 fmadd.d fs4, fa6, fa4, fs4	980 fsub.d ft4, ft4, fs8
···	832 TMadd.d TS3, Tab, Ta3, TS3	884 fmul.d fa6, fa1, fa2	932 TMadd.d TS5, Ta6, Ta5, TS5	981 fmul.d ft5, ft5, ft6
781 TMUL.0 TC3, Td0, TC0	833 TMadd.d fs4, fa6, fa4, fs4	885 fadd.d ft9, ft0, ft1	933 fsgnj.d fa3, ft0, ft0	982 fmsub.d ft10, fa2, ft1, ft0
782 fout d w fs0 zero	834 Tiliduu.u TS5, Td0, Td5, TS5 825 fld fo1 16(+0)	886 TSUD.0 TT10, TT0, TT1	934 (Sgilj.u (S11, (C1, (C1)))	983 fmul.d ft7, ft7, ft6
784 fout d w fs1	836 fld $f_{32} 8(+0)$	887 TSUD.U TIII, TUU, TUI 888 fould d for ft 6 ft 1	955 TSglij.u TS10, TC0, TC0	984 fmul.d ft11, ft6, fs8
$785 \text{ fevt } d \text{ w fs}^2$	837 feed d fe6 ft0 ft1	800 fould d for f = 0	938 fmul d fs9 fa3 fa3	985 fsub.d ft4, ft3, ft4
786	838 fsub d fs7 ft0 ft1	889 ISUD.U 184, ICO, ICI	939 fmul d fa 3 fa 2 fa 3	986 fmadd.d ft8, ft5, ft3, ft8
787 fcvt.d.w fs3. zero	839 fsub.d fs8, ft0, ft1	891 fsub.d fa5 f+0 f+1	940 fmul.d fs6, fs11, fs11	987
788 fcvt.d.w fs4, zero	840 fsub.d fs9, ft0, ft1	892 fmadd.d fs0, fa6, fs6, fs0	941 fmul.d fa4, fa2, fs10	988 fmadd.d ft10, ft7, ft3. ft10
789 fcvt.d.w fs5, zero	841 fsub.d fs10, ft0, ft1	893 fmadd.d fs1, fa6, fs7, fs1		989 addi a0, s0, 920
790	842 fsub.d fs11, ft0, ft1	894 fmadd.d fs2, fa6, fs8, fs2	944 fmul.d fa3, fa3, ft1 🖌	990 scfgwi
791 fsgnj.d ft4, ft3, ft3	843 fmul.d fa6, fa1, fa2	895 fmadd.d fs3, fa6, fs9, fs3	945 fsub.d fs8, fa1, fs9	fnmsub d ft11 ft4 ft6 ft1
792 fsgnj.d ft5, ft3, ft3 🗸	844 fadd.d ft9, ft0, ft1	896 fmadd.d fs4, fa6, fs10, fs4	946 fsub.d fs7, fa1, fs6	001 cofqui
793 fsgnj.d ft7, ft6, ft6	845 fsub.d ft10, ft0, ft1	897 fmadd.d fs5, fa6, fs11, fs5	947 add a1, s10, a5	992 scfaui
794 fsgnj.d ft8, ft6, ft6	846 fsub.d ft11, ft0, ft1	898 fnmsub.d fs0, fa6, ft9, fs0	fmul.d fs6, fs9, fs6	992 Scigwi
795 fld fa1, 8(ra)	847 fsub.d fa3, ft0, ft1	899 fmadd.d fs1, fa6, ft10, fs1	948 addi a0, a1, 912	555 madu.u 1010, 100, 105, 100
796 fadd.d fs6, ft0, ft1	848 fsub.d fa4, ft0, ft1	900 auipc a1, 0x1	fmul.d fa4, fa4, fs11	
797 fadd.d fs7, ft0, ft1	849	fmadd.d fs2, fa6, ft11, fs2	949 add a3, s11, a5	995 Tadd.d Tt8, Tt8, Tt11
798 tadd.d ts8, tt0, tt1	850 fsub.d fa5, ft0, ft1	901 addi a1, a1, -1888	fmul.d fs10, fa3, fs10	
799 fadu.u fs9, ft0, ft1	851 TMadd.d TS0, Tab, TS0, TS0	tmadd.d ts3, ta6, ta3, ts3	950 a001 a7, a3, 912	998 tsa tt10, 0(a0)
801 fadd d fs11 ft0 ft1	852 finadu.u 151, 180, 157, 151 853 fmodd d fs2 fa6 fs8 fs2	902 Tillduu.u TS4, Td0, Td4, TS4	051 fmul d fs7 fs9 fs7	
802 fmadd.d ft3, fs6, fa1, ft3	854 fmadd d fs3 fa6 fs9 fs3	fmadd d fs5 fa6 fa5 fs5	951 mul.u 137, 139, 137	1000 fsd ft8, 0(a/) 🧳
803 fmadd.d ft4, fs7, fa1, ft4	855 fmadd d fs4 fa6 fs10 fs4	904 addi a1 a1 -1892	953 fmul.d fs9, fs9, fa4	
804 fmadd.d ft5, fs8, fa1, ft5	856 fmadd.d fs5, fa6, fs11, fs5	fld = fa7, 16(t0)	954 fmul.d ft9, fs8, ft5	
805 fmadd.d ft6, fs9, fa1, ft6	857 fnmsub.d fs0, fa6, ft9, fs0	905 fld fa1, 0(a1)	955	
806 fmadd.d ft7, fs10, fa1, ft7	858 fmadd.d fs1, fa6, ft10, fs1	906	956 fmul.d ft10, fs7, ft4	
807 fmadd.d ft8, fs11, fa1, ft8	859 fmadd.d fs2, fa6, ft11, fs2	907 fld fa2, 0(a1)	957 fmul.d ft11, fs10, fs1	
808 fld fa1, 8(t0)	860 fmadd.d fs3, fa6, fa3, fs3	908 fsub.d fs6, ft0, ft1	958 fmul.d fs8, fs8, ft8	
809 fsub.d fs6, ft0, ft1	861 fmadd.d fs4, fa6, fa4, fs4	909 fsub.d fs7, ft0, ft1	959 fmul.d fs7, fs7, ft7	
810 fsub.d fs7, ft0, ft1	862 fmadd.d fs5, fa6, fa5, fs5	910 fsub.d fs8, ft0, ft1	960 fmadd.d ft9, fs9, fs0, ft9	
811 fsub.d fs8, ft0, ft1	863 fld fa1, 16(ra)	911 fsub.d fs9, ft0, ft1	961 fmadd.d ft10, fs6, ft3, ft10	
812 tsub.d fs9, ft0, ft1	864 fadd.d fs6, ft0, ft1	912 fsub.d fs10, ft0, ft1	962 fmadd.d ft11, fs11, fs2, ft11	
813 tsub.d ts10, tt0, tt1	865 fadd.d fs7, ft0, ft1	913 fsub.d fs11, ft0, ft1	963 tmadd.d fs8, fs9, fs3, fs8	
814 TSUD.0 TS11, Tt0, Tt1	866 tadd.d ts8, tt0, tt1	914 tmul.d ta6, fa7, fa7	964 tmadd.d ts/, ts6, tt6, fs7	
ALS TIMULA TAB, TAL, TAL	86/ tadd.d ts9, tt0, tt1	915 tadd.d tt9, tt0, tt1	965 tmul.d ts0, ts10, ts4	
OLD TAUL.U TTY, TTU, TTL 217 fould d ft10 ft0 ft1	808 Taud.d TS10, TT0, TT1	916 tsub.d tt10, tt0, tt1	900 Taud.d Tt9, Tt9, Tt10	
$\frac{017}{1500.0}$ TL10, TL0, TL1 818 feath d ft11 ft0 ft1	70 Taula. TSII, TTU, TTI 870 fmodd d ft2 fc6 fo1 ft2	91/ TSUD.0 TTII, TTU, TTI 018 fould for ft0 ft1	907 Tauu.u Tt3, Tt3, Tt4 968 fadd d ft4 ft6 ft7	
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	871 fmodd d ft4 fc7 fo1 ft4	910 TSUD.0 Td3, TT0, TT1 010 fould d fo4 ft0 ft1	500 radu.u TL4, TL0, TL7	
820 fsub.d fa4. ft0. ft1	872 fmadd.d ft5, fs8, fa1 ft5	920 fsub d fa5 ft0 ft1	707 Iduu.u TSO, TSO, TS/	
821	873 fmadd.d ft6, fs9, fa1, ft6	921 fmadd.d fs0. fa6. fs6 fs0	970 THADA.A TSU, TSII, TS5, TSU	
021	5.5	130 jiao jiao jiao jiao jiao jiao jiao jiao	9/1 tadd.d tt9, tt9, tt11	

Add *date* or a third information here

PsPIN: A PULP-powered implementation of sPIN





Stochastic Rounding

 Stochastic rounding to avoid stagnation when accumulating small quantities into a large accumulator

 Rounding decision taken comparing p_rsr bits against the output of an LFSR

When using p_rsr = 12 bits, the area overhead amounts to ~4% on the ExSdotp unit, ~1% on the full extended FPU.



