



Thousand-Way Parallel RTL Simulation



Mahyar Emami, Thomas Bourgeat, James R. Larus

Simulation 1S Slow







3

A Growing Gap





4

A Growing Gap







1 GHz chip at 1000 Hz \rightarrow million times slowdown





1 GHz chip at 1000 Hz \rightarrow million times slowdown





1 GHz chip at 1000 Hz \rightarrow million times slowdown





Reproduced from https://www.karlrupp.net/2018/02/42-years-of-microprocessor-trend-data/

We need 1000-way parallel simulators!



Parendi: Thousand-Way Parallel RTL Simulation

 \bullet Runs RTL simulation on a real thousand-core chip \rightarrow Graphcore IPU





Parendi: Thousand-Way Parallel RTL Simulation

- Runs RTL simulation on a real thousand-core chip \rightarrow Graphcore IPU
- Graphcore IPU (Intelligence Processing Unit)
 - 1472 cores per socket
 - \circ ≈ 600 KiB of SRAM per core \rightarrow 900 MiB on-chip
 - Message-passing architecture
 - Bulk-synchronous parallel (BSP) execution model
 - Low-cost synchronization/communication
 - Scale out to 4x (this work), 16x, 64x, ...









Parendi: Thousand-Way Parallel RTL Simulation

- Runs RTL simulation on a real thousand-core chip \rightarrow Graphcore IPU
- Graphcore IPU (Intelligence Processing Unit)
 - 1472 cores per socket
 - \circ ~ 600 KiB of SRAM per core \rightarrow 900 MiB on-chip
 - Message-passing architecture
 - Bulk-synchronous parallel (BSP) execution model
 - Low-cost synchronization/communication
 - \circ Scale out to **4x (this work)**, 16x, 64x, ...

• Parendi

- Targets the IPU (based on Verilator)
- 1000-way parallelization and optimizations
- Up to 4x faster than Verilator on x64









Background





Simulating Chips: What is RTL?

- A netlist \rightarrow a cyclic graph
- "Stateless" combinational logic
- "Stateful" registers and memories
 - Retain values at clock cycle 0



Digital hardware is described in RTL (register transfer level)





Cycle-accurate RTL Simulation*

- Split registers
 - **current** and **next** value



Netlist Graph

* Not timing accurate, but good enough for almost all functional verification tasks.



Cycle-accurate RTL Simulation*

- Split registers
 - **current** and **next** value



* Not timing accurate, but good enough for almost all functional verification tasks.

At each cycle:

- Compute next
- Overwrite current
- Repeat





Computation phase

- Independently compute program partitions
- Globally synchronize





Computation phase

- Independently compute program partitions
- Globally synchronize







Computation phase

- Independently compute program partitions
- Globally synchronize



Communication phase

- Exchange produced values
- Globally synchronize





Computation phase

- Independently compute program partitions
- Globally synchronize



How does this scale to 1000s of cores/threads?

Communication phase

- Exchange produced values
- Globally synchronize

Thousand -way Parallel Simulation

Scaling Simulation Performance number of cores r(N) =simulation rate Minimize Comp. Syn Comm.

Thread 1 a3 b3

Thread 2 a4 a5 b4

RTL Cycle 1

Scaling Simulation Performance number of cores r(N)simulation rate Comp. Syn Comm. Thread 0 al a2 a3 b1 b2 c1 $R_{0-} = R_{0+}$ a1 a2 a3 b1 b2 c1 • • • Thread 1 a3 b3

Thread 2 a4 a5 b4

Why use the IPU? What about a general purpose processor?

 $t_{comp} + t_{sync} + t_{comm}$ Minimize

23

IPU or x64

	computation	synchronization	communication
IPU	1000s of slow cores X	≃const ✔	On-chip ≃ const ✔ Off-chip varies X
x64	10s/100s of fast cores 🖌	increasing with N X	increases with N X

$$\frac{1}{_{mp}+t_{sync}+t_{comm}}$$

IPU or x64

	computation	synchronization	communication
IPU	1000s of slow cores X	≃const ✔	On-chip ≃ const ✔ Off-chip varies X
x64	10s/100s of fast cores 🖌	increasing with N X	increases with N X

See paper for the quantitative treatment

$$\frac{1}{_{mp}+t_{sync}+t_{comm}}$$

Overview of Parendi

DAG

Build a data dep. graph from source code (Verilog) 1.

Tiles

Overview of Parendi

DAG

Fibers

- Build a data dep. graph from source code (Verilog) 1.
- 2.

Tiles

Extract *fibers* \rightarrow smallest units of independent work

Overview of Parendi

Fibers

- Build a data dep. graph from source code (Verilog) 1.
- 2.
- 3.

Extract *fibers* \rightarrow smallest units of independent work Partition fibers across IPUs → hypergraph partitioning

Compiling RTL to the IPU

Fibers

4.	Partition fibers within]	CF
3.	Partition fibers across]	EF
2.	Extract <i>fibers</i> \rightarrow smalles	;t
1.	Build a data dep. graph f	Fr

IPUs

Tiles

rom source code (Verilog) : units of independent work PUs \rightarrow hypergraph partitioning PUs \rightarrow submodular load balancing

Evaluation

Evaluation

- $\{1|s\}rN \rightarrow RV64 \text{ rocket}^1 \text{ cores on }NxN$ mesh²
- $lr \rightarrow Large rocket cores$
- sr \rightarrow Small rocket cores
- Light-weight Verilog harness

 [1] Krste Asanović et al, "The Rocket Chip Generator." (2016).
 [2] Jerry Zhao et al, "Constellation: An Open-Source SoC-Capable NoC Generator," 2022 15th IEEE/ACM International Workshop on Network on Chip Architectures (NoCArc 2022).

Hardware Setup

	Verilator (v5.006)		Parendi
Hardware	AMD EPYC 9954 (ae4)	Intel Xeon 6348 (ix3)	4xIPU M2000 (ipu)
<pre># cores avail.</pre>	64x2	28x2	1472x4
# cores used	1, 2, 4, 6,, 32	1, 2, 4, 6,, 32	1472, 2944, 4416, 5888
Freq. GHz	3.75	3.5	1.35
SRAM (MiB)	386x2	79x2	900x4
Released	Q4 2022	Q2 2021	Q3 2020

Verilator parallelism limited to 32 cores due to excessive compile time (+8hr); see paper.

Conclusion

- Thousand-way parallel RTL simulation is becoming a necessity
- Parendi shows thousand-way parallelism is feasible
- Runs on real hardware \rightarrow Graphcore IPU
- 5888 cores across 4 sockets on a server board
- Multi-stage parallelization algorithm
- Up to 4x faster than Verilator on high-end x64
 - See paper for more evaluation, including cost, other academic work, microbenchmarks and etc.

Backup

Cost Analysis

- Regression testing scenario
- N identical tests \circ short \rightarrow 1 million cycles
- Ad-hoc:
 - IPU: One test per IPU (1472 tiles)
 - x64: One test per CPU core
- Fine:
 - Tests parallelized to multiple
 optimal (perf.) tiles and cores

Weak Scaling

 $r(N) = \frac{1}{t_{co}}$

- Increase total work and paral
 Computation time remains cor
 Simplification → no communic
- \bullet Weak scaling goal \rightarrow constant sim. rate w.r.t to N

$$\frac{1}{t_{mp} + t_{sync} + t_{comm}}$$
I elism proportionally istant w.r.t to N ation cost

Sync. Cost Analysis

Sync. Cost Analysis

Sync. Cost Analysis

Sync. Cost Analysis

The IPU has far more scalability potentials: sync. cost is not an issue!

