# BOOM
# An open-source out-of-order processor

Christopher Celio, Jerry Zhao, Abraham Gonzalez,
Ben Korpan, Krste Asanovic, David Patterson

## Chisel Community Conference 2018

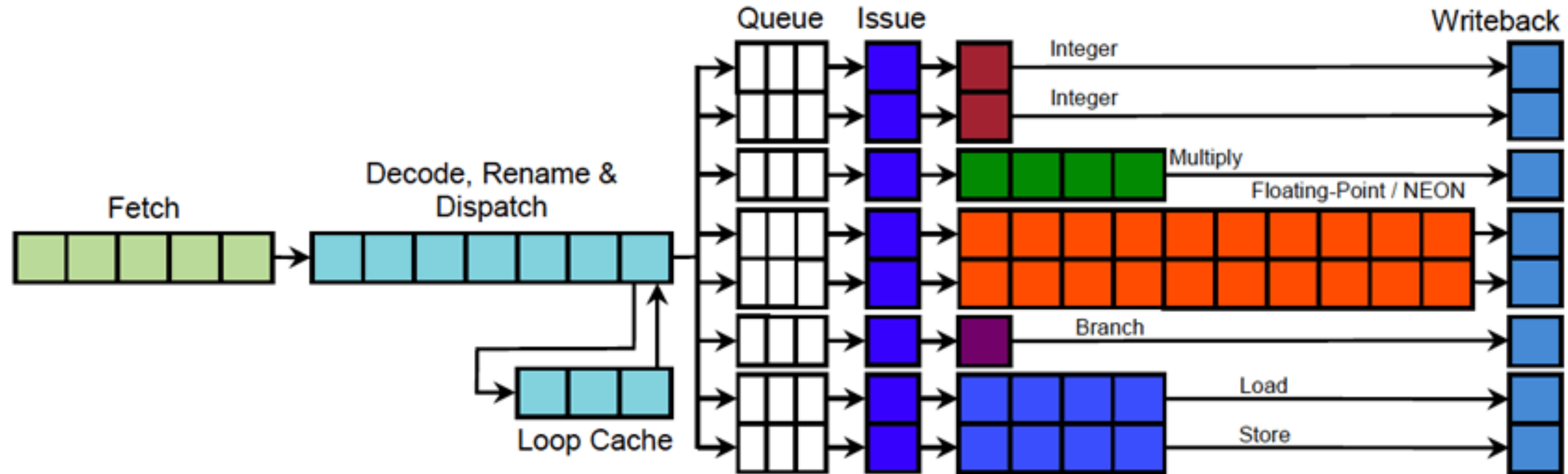https://github.com/riscv-boom

# Outline

- Motivation

- Microarchitectural Overview

- Current state of the project

- How to get started

- How to contribute

- Things to work on

- Discussion
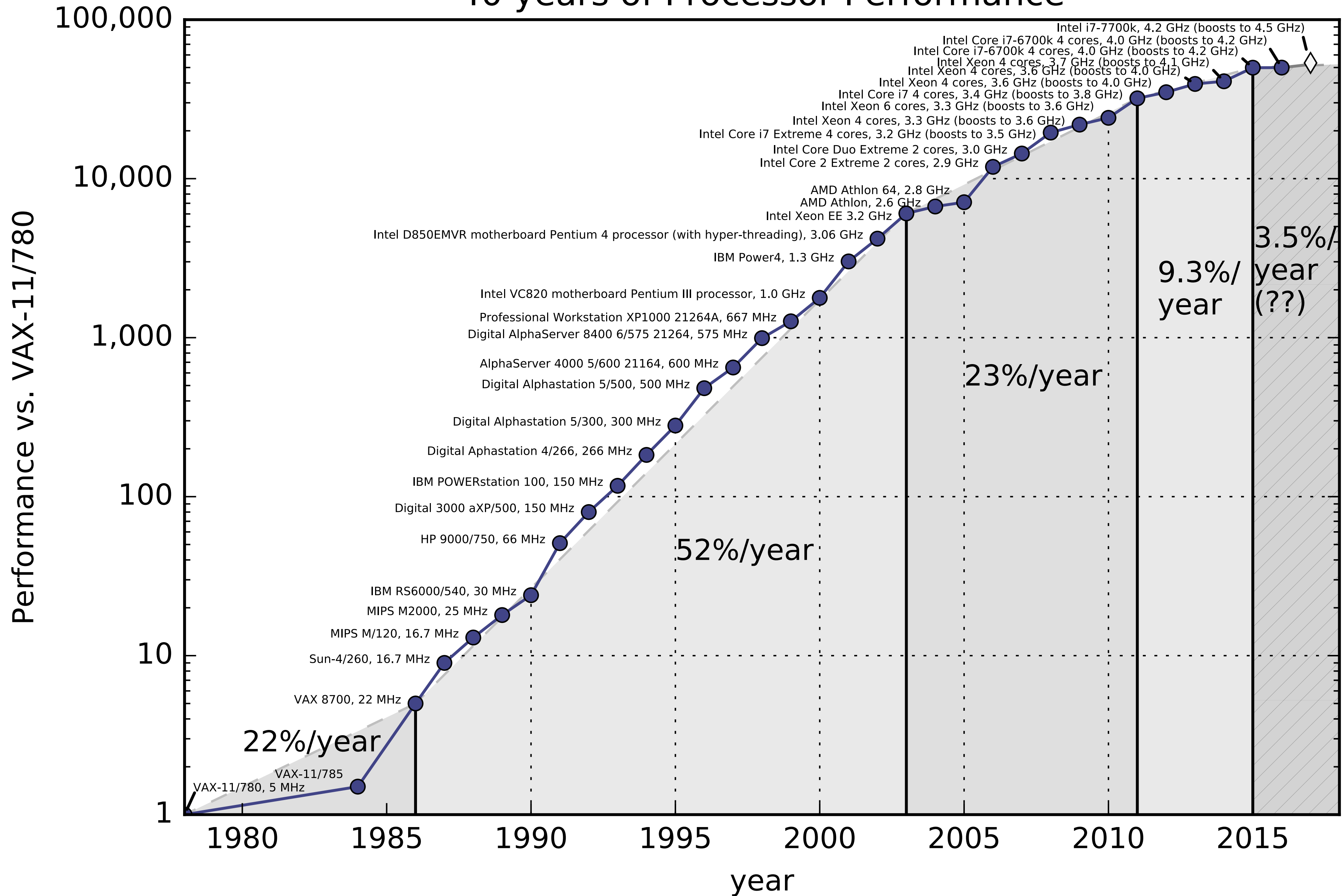
# What does a modern processor look like?



**ARM Cortex A-15**

- pipelining
  - overlap execution of multiple instructions
- superscalar
  - multiple instructions / cycle / stage
- out-of-order
  - execute instructions in dependence order

# 40 years of Processor Performance

**Performance vs. VAX-11/780** (y-axis)

**year** (x-axis)

- Intel i7-7700k, 4.2 GHz (boosts to 4.5 GHz)
- Intel Core i7-6700k 4 cores, 4.0 GHz (boosts to 4.2 GHz)
- Intel Core i7-6700k 4 cores, 4.0 GHz (boosts to 4.2 GHz)
- Intel Xeon 4 cores, 3.7 GHz (boosts to 4.1 GHz)
- Intel Xeon 4 cores, 3.6 GHz (boosts to 4.0 GHz)
- Intel Xeon 4 cores, 3.6 GHz (boosts to 4.0 GHz)
- Intel Core i7 4 cores, 3.4 GHz (boosts to 3.8 GHz)
- Intel Xeon 6 cores, 3.3 GHz (boosts to 3.6 GHz)
- Intel Xeon 4 cores, 3.3 GHz (boosts to 3.6 GHz)
- Intel Core i7 Extreme 4 cores, 3.2 GHz (boosts to 3.5 GHz)
- Intel Core Duo Extreme 2 cores, 3.0 GHz
- Intel Core 2 Extreme 2 cores, 2.9 GHz
- AMD Athlon 64, 2.8 GHz
- AMD Athlon, 2.6 GHz
- Intel Xeon EE 3.2 GHz
- Intel D850EMVR motherboard Pentium 4 processor (with hyper-threading), 3.06 GHz
- IBM Power4, 1.3 GHz
- Intel VC820 motherboard Pentium III processor, 1.0 GHz
- Professional Workstation XP1000 21264A, 667 MHz
- Digital AlphaServer 8400 6/575 21264, 575 MHz
- AlphaServer 4000 5/600 21164, 600 MHz
- Digital Alphastation 5/500, 500 MHz
- Digital Alphastation 5/300, 300 MHz
- Digital Alphastation 4/266, 266 MHz
- IBM POWERstation 100, 150 MHz
- Digital 3000 aXP/500, 150 MHz
- HP 9000/750, 66 MHz
- IBM RS6000/540, 30 MHz
- MIPS M2000, 25 MHz
- MIPS M/120, 16.7 MHz
- Sun-4/260, 16.7 MHz
- VAX 8700, 22 MHz
- VAX-11/785
- VAX-11/780, 5 MHz

22%/year

52%/year

23%/year

9.3%/year

3.5%/year (??)

Data (mostly) from H&P CAAQA 6th Ed

**4**
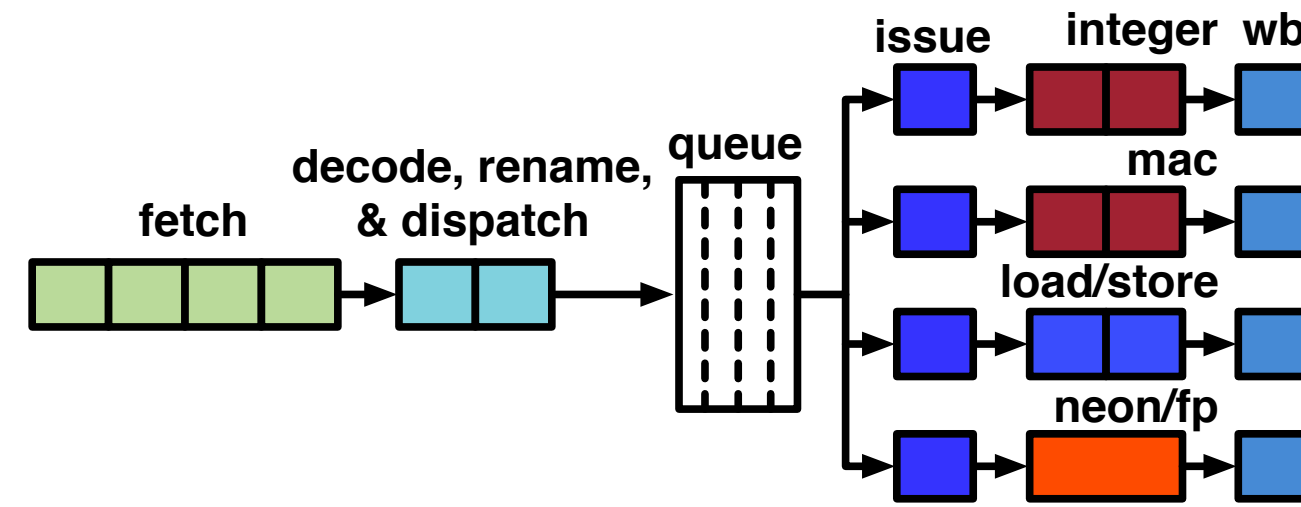
# Academic OoO Research

- Lack of effort in academia to build, evaluate OoO designs
  - NCSU's Fabscalar/AnyCore
  - MIT's riscy-OOO (bluespec)
- most research uses software simulators
  - SimpleScalar (5000 cites), GEM5 (1000 cites), SESC (250 cites)
  - generally don't support full systems
  - cannot produce area, power numbers
  - hard to trust, verify results
    - McPAT is calibrated against 90nm Niagara, 65nm Niagara 2, 65nm Xeon, and 180nm Alpha 21364
  - very slow (~100 KHz)
    - ~10 hours of sim is 1 second of target
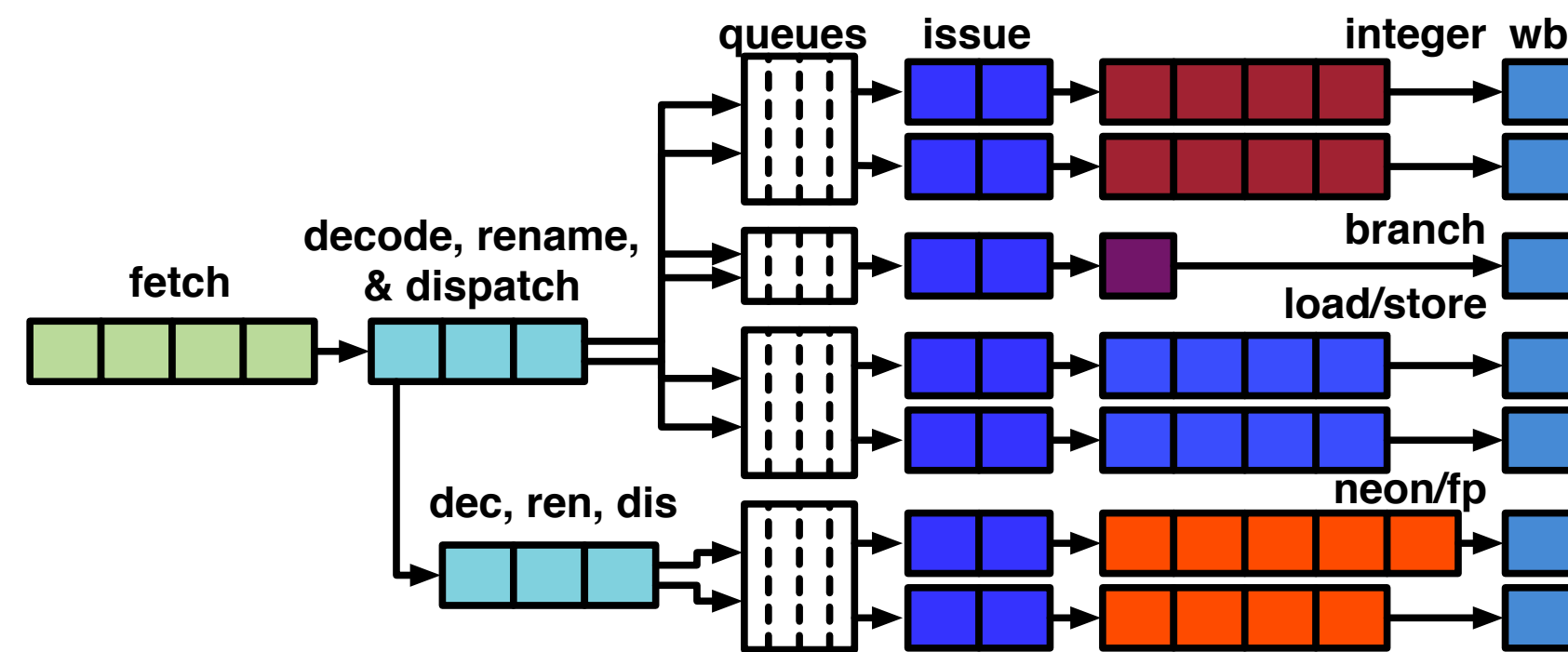
**2007**



**Cortex-A9**

**2012**



**Cortex-A15**

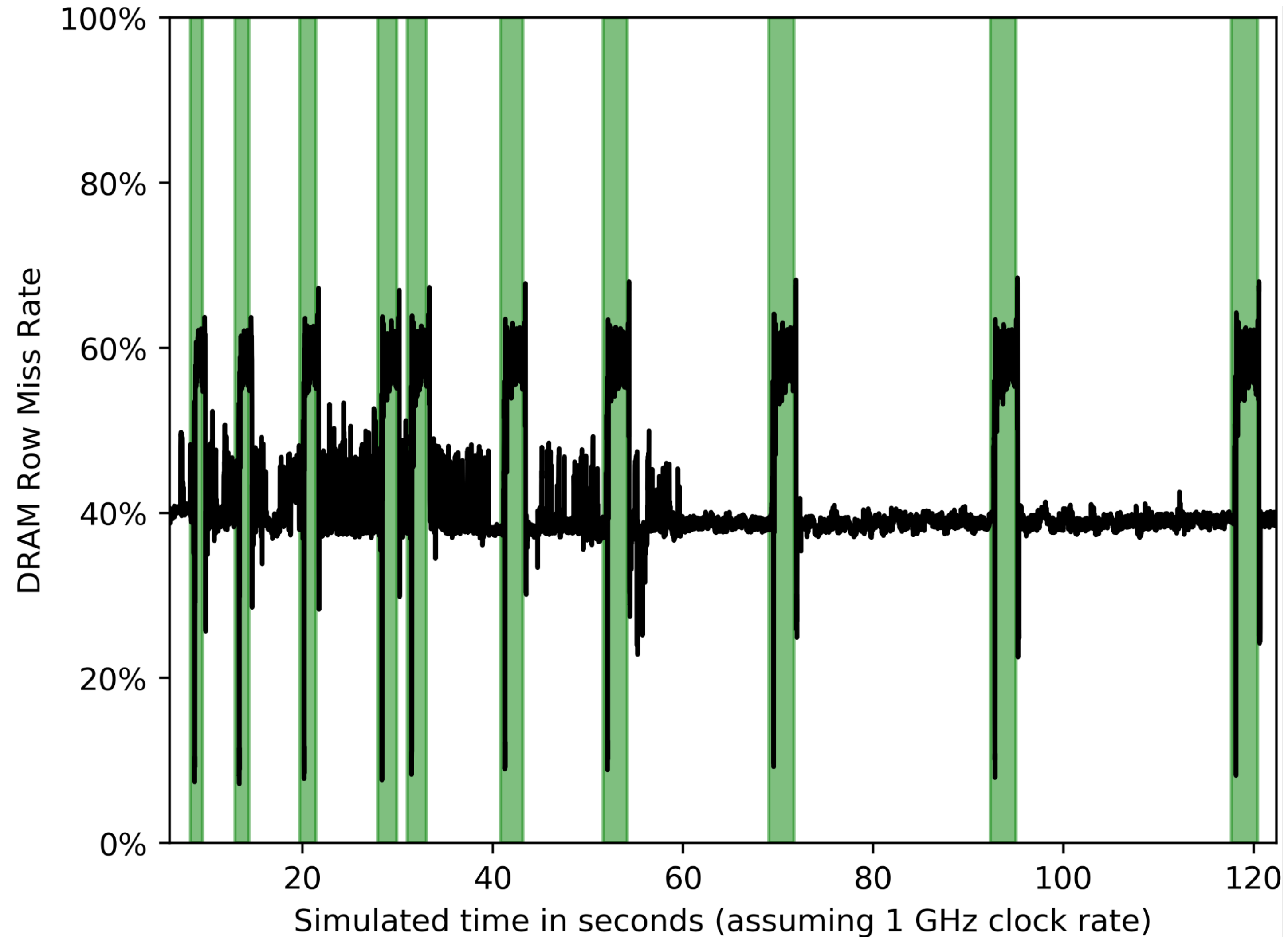**2016**



**Cortex-A73**

# Berkeley Architecture Research

- RTL of processor system
- FPGA-based simulators (50-100 MHz)
- **trillions** of instructions simulated (full workloads)
- power models built from actual activity
- can generate floor-plans, area, timing reports

data collected by
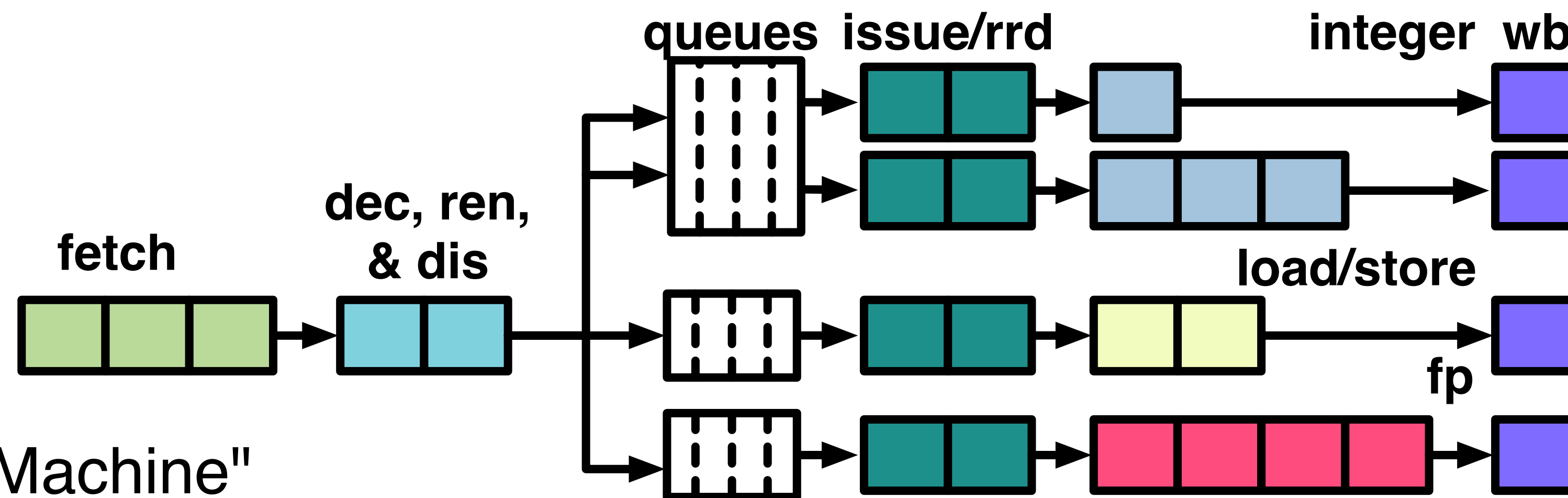Martin Maas, et. al.

# What is BOOM useful for?

- As a complex IP for methodology studies
  - How do you model RTL performance?
  - How do you measure RTL power?
  - How can we help build an agile verification story?
  - How do we do open-source cad flows?
- Software Studies
  - hardware/software co-design
  - high visibility of very long-running applications
- Off-the-shelf out-of-order core
  - need a core for your research tape-out to talk to your IP?
  - drive memory-mapped accelerators
- Realistic Microarchitecture Studies of contained blocks
  - branch prediction, issue queue design, etc.
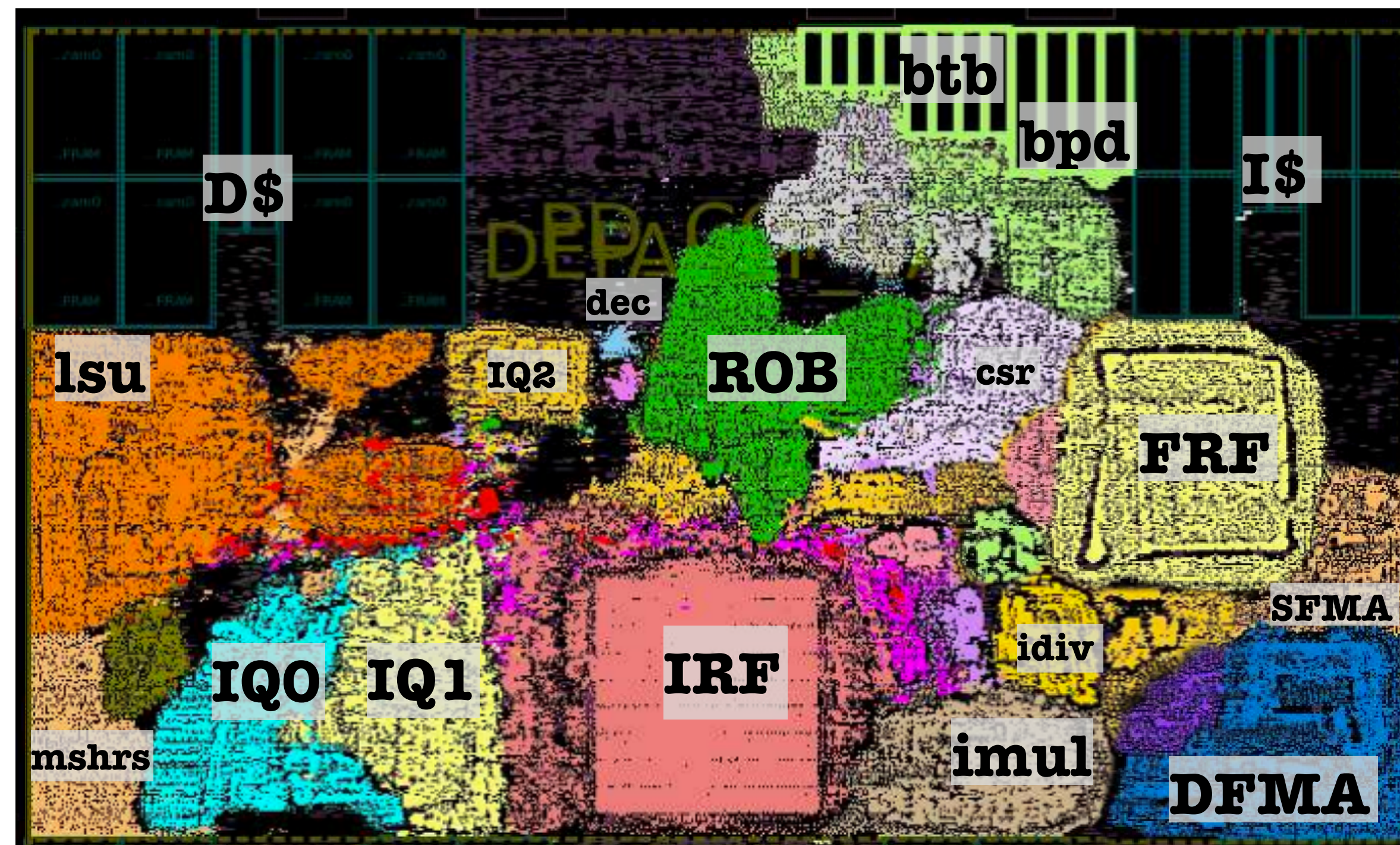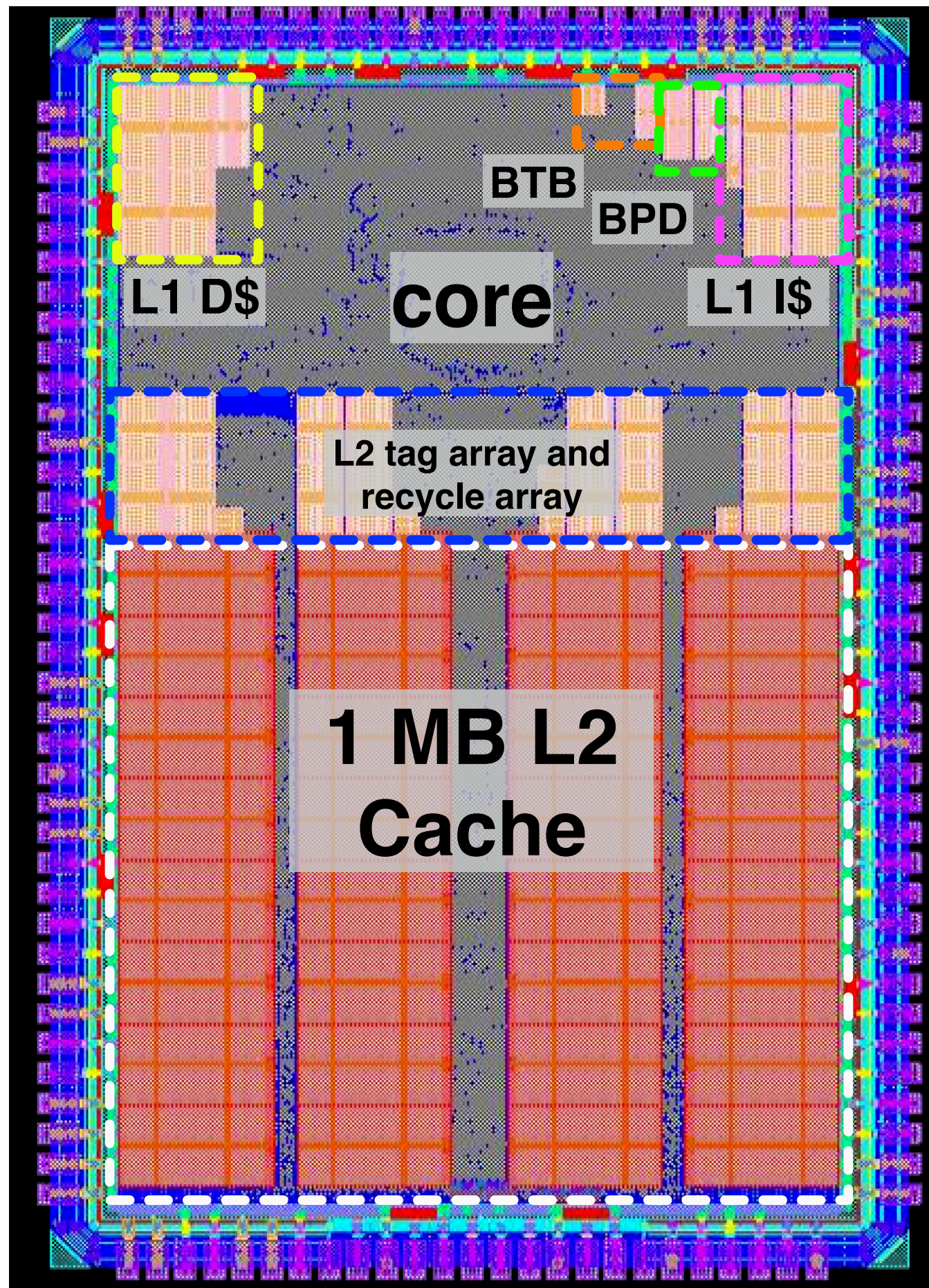- Security Research

# What is BOOM?



- "Berkeley Out-of-Order Machine"
- superscalar
- out-of-order
- implements **RV64G**, boots Linux
- It is synthesizable
- it is open-source
- written in **Chisel** (16k loc)
- It is parameterizable generator
- built on top of Rocket-chip SoC Ecosystem
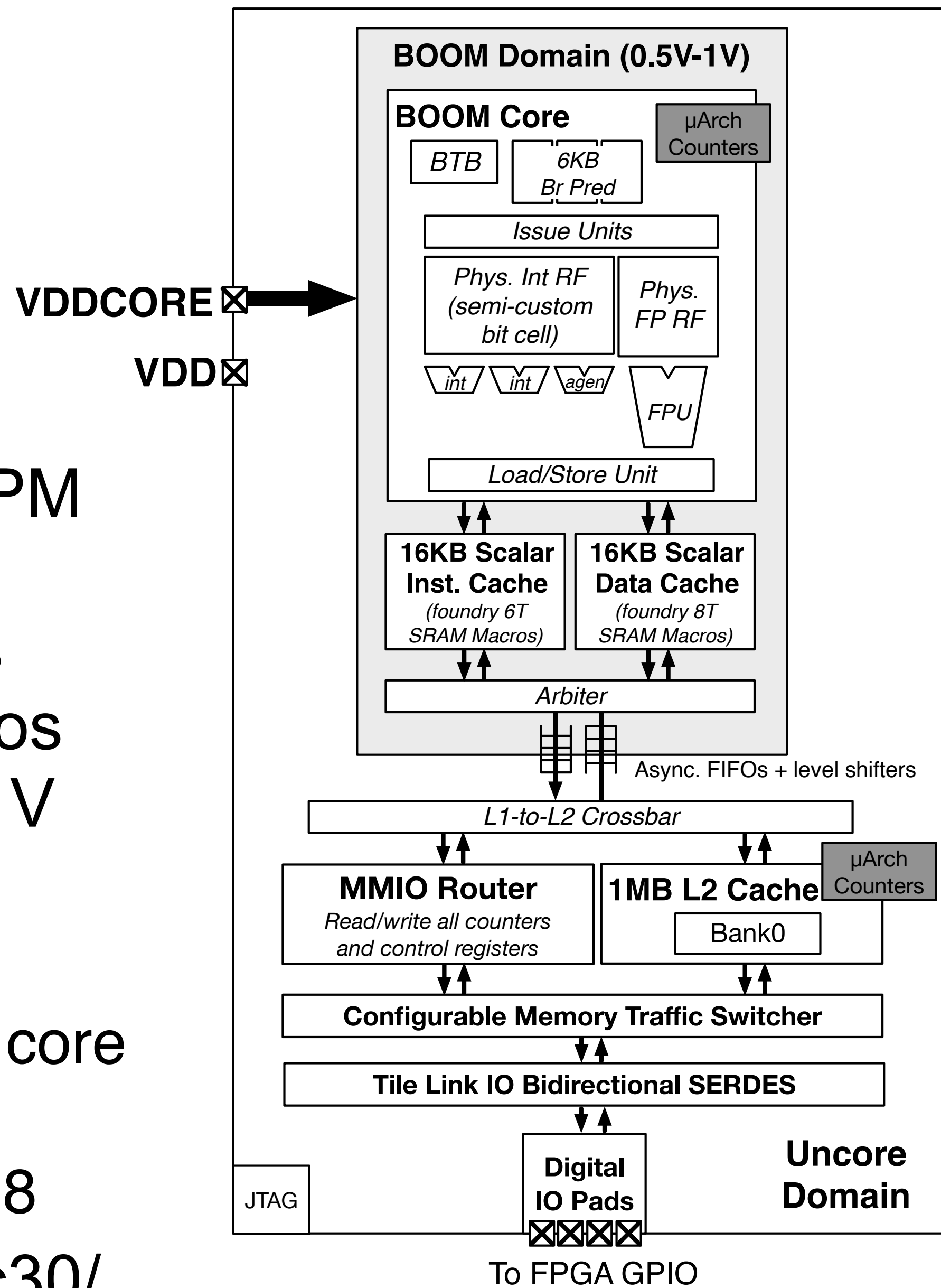
http://ucb-bar.github.io/riscv-boom

# BROOM Chip (Taped out Aug 2017)



TSMC 28 nm HPM
$6\ mm^2$
417k std cells
73 SRAM macros
1.0 GHz @ 0.9 V

- Open-source superscalar out-of-order RISC-V core
- Resilient cache for low-voltage operation
- See BROOM work presented in Hot Chips 2018

https://www.hotchips.org/archives/2010s/hc30/

# **Leveraging Open-source RTL**

- The Rocket-chip SoC Generator
- Started in 2011
- Taped out ~~10~~ (~~13?~~) *17* times by Berkeley + many others
- 6,257 commits
- 72 contributors
- Commercial quality
- Replace standard in-order core with BOOM
- Leverage Rocket-chip as a library of processor components



**BOOM goes here**

https://github.com/freechipsproject/rocket-chip

# The BOOMv2 Core

# Core Comparison

| Processor | SiFive U54 Rocket (RV64GC) | Berkeley BOOMv2 (RV64G) | OpenSPARC T2 | ARM Cortex-A9 | Intel Xeon Ivy |
|---|---|---|---|---|---|
| Language | Chisel | Chisel | Verilog | - | SystemVerilog |
| Core LoC | 8,000 | 16,000 | 290,000 | - | - |
| SoC LoC | 34,000 | 50,000 | 1,300,000 | - | - |
| **Foundry** | TSMC | TSMC | TI | TSMC | Intel |
| **Technology** | 28 nm (HPC) | 28 nm (HPM) | 65 nm | 40 nm (G) | 22 nm |
| **Core+L1 Area** | 0.54 mm$^2$ | 0.52 mm$^2$ 16kB/16kB | ~12 mm$^2$ | ~2.5 mm$^2$ | ~12 mm$^2$ core+L1+L2 |
| **Coremark/MHz** | 2.75 | 3.77 | 1.64* | 3.71 | 5.60 |
| **Frequency** | 1.5 GHz | 1.0 GHz | 1.17 GHz | 1.4 GHz | 3.3 GHz |

*From eembc.org. 32 threads/8 cores achieve 13 Cm/MHz.

14

# The Evolution of BOOM

| | BOOMv1 | BOOMv2 |
|---|---|---|
| BTB entries | 40 (fully-assoc) | 64 x 4 (set-assoc) |
| Fetch Width | 2 insts | 2 insts |
| Issue Width | 3 micro-ops | 4 micro-ops |
| Issue Entries | 20 | 16/20/10 |
| Regfile | 7r3w (unified) | 6r3w (inst), 3r2w (fp) |
| Exe Units | iALU+iMul+FMA iALU+fDiv Load/Store | iALU+iMul+iDiv iALU FMA+fDiv Load/Store |



BOOM v1 (April 2017)

BOOM v2 (Aug 2017)

# BOOMv3 (tentative)

- Goal
  - Use lessons learned from BOOMv2 tape-out to improve core.
  - Update to the latest RISC-V standards.
- Privileged Spec v1.11, User 2.3
  - Done.
- RISC-V Compressed support
  - TODO.
- RISC-V WMO memory consistency model
  - Must order loads to the same address.
  - TODO.
- updated front-end and branch prediction
  - Performance debugging needed.
- 4-cycle load-use
  - speculates load-hit to save two cycles.
  - Done.

# Current frontend



- **BTB (branch target buffer)**
  - predicts without seeing instructions
  - set-associative, partially tagged
  - checker to verify integrity
- **BIM (bimodal)**
  - a table of two-bit counters
  - used by BTB and optionally BPD to decide direction
- **RAS (return address stack)**
  - predicts returns
  - driven by BTB to make decisions

# Current frontend



- **BPD (conditional predictor)**
  - provide your own (e.g., gshare or TAGE)
  - decides taken/not-taken based on instruction bits
  - uses path history
- **Fetch Target Queue**
  - stores fetch PC, branch prediction information
  - one entry == one fetch bundle

# Abstract Branch Predictors

*holds both rename-table snapshots, and bpd snapshot information that can be released once a branch is resolved.*

*branch snapshots*

*B-ROB holds all inflight branches, and can bypass updates that haven't been committed to incoming predictions.*

```
global history register
* update at end of BP1 stage
   – it contains the branch prediction history
      of the fetch unit (including BTB
      decisions)
   – compresses entire fetch–packet decision
   – only includes branches, not JAL/JALRs
* reset on fetch unit redirect (misprediction)
   – new prediction must use the new history
```

**19**

# GShare in single-ported SRAM

**_Execution Pipeline_**

**_Issue_**   **_RegisterRead_**   **_Execute_**   **_Memory_**   **_Sign-extend_**   **_Writeback_**



- Speculatively wakeup uops that depend on loads
- Kill them on the next cycle if miss occurs

- FPU needs 3 sources
- Only support one Mem unit (one load/store)
- ALU is padded out to max latency
- div unit is unpipelined, can apply back-pressure



Quad-issue (8r,4w)

*bypassing*

ALU

*bypassing*

ALU

FPU

imul

*bypassing*

ALU

div

Agen    LSU    D$

Issue Select

Regfile Read (8 Read Ports) bypassable

bypass network

Regfile Writeback (4 Write Ports)

# Branches

- MIPS R10K style
- Every branch:
  - is given a tag
  - takes a snapshot of the rename map tables
  - is given an empty "allocation list"
- Following instructions:
  - all uops have a branch-mask
  - if bit is set, they depend on that unresolved branch
- When branch is resolved:
  - the branch tag is broadcast across the machine
  - everyone clears their bit
  - allows new allocation
- When branch is mispredicted:
  - all uops with matching branch tag are killed immediately
  - rename tables are set to the snapshot
  - allocated physical registers are added back to free list

# Parameterized Superscalar

**dual-issue (5r,3w)**



```scala
val exe_units = ArrayBuffer[ExecutionUnit]()
exe_units += Module(new ALUExeUnit(is_branch_unit    = true
                   , has_fpu       = true
                   , has_mul       = true
                   ))
exe_units += Module(new ALUMemExeUnit(fp_mem_support = true
                   , has_div       = true
                   ))
```

**Quad-issue (9r,4w)**



## OR

```scala
exe_units += Module(new ALUExeUnit(is_branch_unit = true))
exe_units += Module(new ALUExeUnit(has_fpu = true
           , has_mul = true
           ))
exe_units += Module(new ALUExeUnit(has_div = true))
exe_units += Module(new MemExeUnit())
```

# A Functional Unit



kill
(flush_pipeline)

Br Logic

branch resolution info

req.valid

req.uop.brmask

Op1 Data

Op2 Data

function code

req.ready

FIFOIO

Res → Res → Res → Res → Res

Speculative Pipeline

FIFOIO

resp.valid

WB Data

resp.ready

**Parameters**
num_stages = 4
is_var_latency = false
is_pipelined = true
shares_writeport= false

earliest_bypass_stage=4

is_branch_unit = true

(val, pdst, data)

(val, pdst, data)

(val, pdst, data)

(val, pdst, data)

**BYPASSES**

**25**

- **Abstract FunctionalUnit**
  - describes common IO
- **Pipelined/Iterative**
  - handles storing uop metadata, branch resolution, branch kills
- **Concrete Subclasses**
  - instantiates the actual expert-written FU
  - no modifications required to get FU working with speculative OoO
  - allows easy "stealing" of external code

**OMG**

(300+ in this file
2061 lines more)

```
// See LICENSE for license details.

//*** THIS MODULE HAS NOT BEEN FULLY OPTIMIZED.
//*** DO THIS ANOTHER WAY?

package hardfloat

import Chisel._
import Node._
import consts._

object MaskOnes
{
  def apply(in: UInt, start: Int, length: Int): UInt = {
    val top = 1 << in.getWidth
    val shift = SInt(BigInt(-1) << top) >> in
    Reverse(shift(top-1-start,top-length-start))
  }
}

object estNormDistPNNegSumS
{
  def priorityEncode(key: UInt, n: Int, s: Int) = {
    if (Module.backend.isInstanceOf[CppBackend]) UInt(n+s-1) -
Log2(key(s-1,0), s)
    else PriorityMux((0 until s).map(i => (key(s-1-i), UInt(n+i,
log2Up(n+s-1)))))
  }

  def apply(a: UInt, b: UInt, n: Int, s: Int) =
    priorityEncode((a ^ b) ^ ~((a & b) << UInt(1)), n, s)
}

object estNormDistPNPosSumS
{
  def apply(a: UInt, b: UInt, n: Int, s: Int) =
    estNormDistPNNegSumS.priorityEncode((a ^ b) ^ ((a | b) << UInt(1)), n, s)
}

class mulAddSubRecodedFloatN_io(sigWidth: Int, expWidth: Int) extends Bundle
{
  val op = UInt(INPUT, 2)
  val a = UInt(INPUT, expWidth+sigWidth+1)
  val b = UInt(INPUT, expWidth+sigWidth+1)
  val c = UInt(INPUT, expWidth+sigWidth+1)
  val roundingMode = UInt(INPUT, 2)
  val out = UInt(OUTPUT, expWidth+sigWidth+1)
  val exceptionFlags = UInt(OUTPUT, 5)
}

class mulAddSubRecodedFloatN(sigWidth: Int, expWidth: Int, speed: Boolean =
false) extends Module {
  val io = new mulAddSubRecodedFloatN_io(sigWidth, expWidth)

  val sigSumSize = (sigWidth+2)*3
  val normSize = (sigWidth+2)*2
  val logNormSize = log2Up(normSize)
  val firstNormUnit = 1 << logNormSize-2
  val minNormExp = (1 << expWidth-2) + 2
  val minExp = minNormExp - sigWidth

  val signA  = io.a(expWidth+sigWidth)
  val expA   = io.a(expWidth+sigWidth-1, sigWidth)
  val fractA = io.a(sigWidth-1, 0)
  val isZeroA = expA(expWidth-1, expWidth-3) === UInt(0)
  val isSpecialA = expA(expWidth-1, expWidth-2) === UInt(3)
  val isInfA = isSpecialA && !expA(expWidth-3)
```
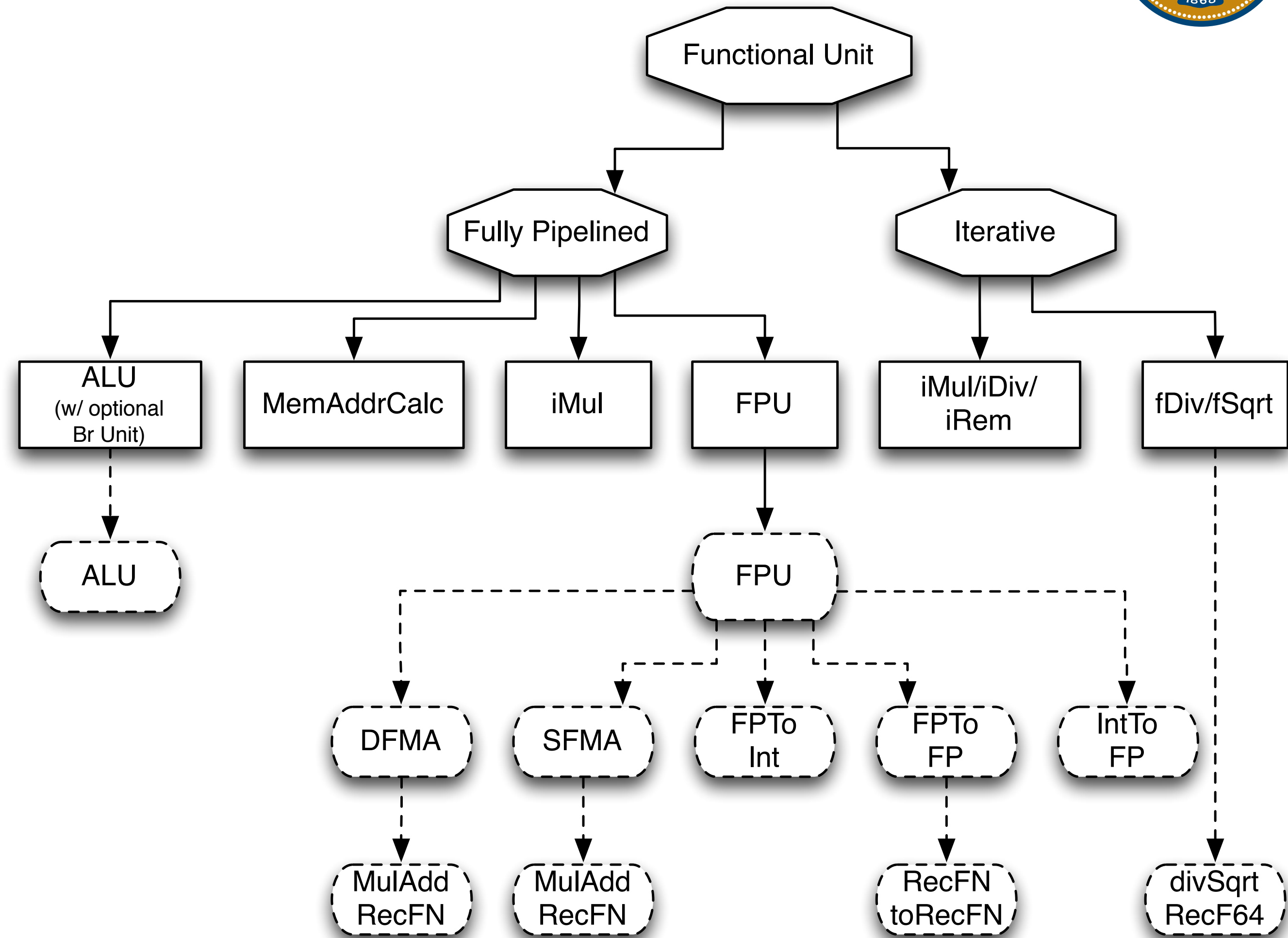
```
val sigA = Cat(!isZeroA, fractA)

val signB  = io.b(expWidth+sigWidth)
val expB   = io.b(expWidth+sigWidth-1, sigWidth)
val fractB = io.b(sigWidth-1, 0)
val isZeroB = expB(expWidth-1, expWidth-3) === UInt(0)
val isSpecialB = expB(expWidth-1, expWidth-2) === UInt(3)
val isInfB = isSpecialB && !expB(expWidth-3)
val isNaNB = isSpecialB && expB(expWidth-3)
val isSigNaNB = isNaNB && !fractB(sigWidth-1)
val sigB = Cat(!isZeroB, fractB)

val opSignC  = io.c(expWidth+sigWidth) ^ io.op(0)
val expC   = io.c(expWidth+sigWidth-1, sigWidth)
val fractC = io.c(sigWidth-1, 0)
val isZeroC = expC(expWidth-1, expWidth-3) === UInt(0)
val isSpecialC = expC(expWidth-1, expWidth-2) === UInt(3)
val isInfC = isSpecialC & !expC(expWidth-3)
val isNaNC = isSpecialC &  expC(expWidth-3)
val isSigNaNC = isNaNC & !fractC(sigWidth-1)
val sigC = Cat(!isZeroC, fractC)

val roundingMode_nearest_even = io.roundingMode === round_nearest_even
val roundingMode_minMag       = io.roundingMode === round_minMag
val roundingMode_min          = io.roundingMode === round_min
val roundingMode_max          = io.roundingMode === round_max

//-------------------------------------------------------------------
//-------------------------------------------------------------------
val signProd = signA ^ signB ^ io.op(1)
val isZeroProd = isZeroA || isZeroB
val sExpAlignedProd = Cat(Fill(3, !expB(expWidth-1)), expB(expWidth-2, 0)) + expA +
UInt(sigWidth+4)

//-------------------------------------------------------------------
//-------------------------------------------------------------------
val doSubMags = signProd ^ opSignC

val sNatCAlignDist = sExpAlignedProd - expC
val CAlignDist_floor = isZeroProd || sNatCAlignDist(expWidth+1)
val CAlignDist_0 = CAlignDist_floor || sNatCAlignDist(expWidth, 0) === UInt(0)
val isCDominant = !isZeroC && (CAlignDist_floor || sNatCAlignDist(expWidth, 0) < UInt(sigWidth+2))
val CAlignDist =
    Mux(CAlignDist_floor, UInt(0),
    Mux(sNatCAlignDist(expWidth, 0) < UInt(sigSumSize-1), sNatCAlignDist,
    UInt(sigSumSize-1)))(log2Up(sigSumSize)-1, 0)
val sExpSum = Mux(CAlignDist_floor, expC, sExpAlignedProd)

// *** USE `sNatCAlignDist'?
var CExtraMask = MaskOnes(CAlignDist, normSize, sigWidth+1)
val negSigC = Mux(doSubMags, ~sigC, sigC)
val alignedNegSigC =
    Cat(Cat(doSubMags, negSigC, Fill(normSize, doSubMags)).toSInt >> CAlignDist,
    ((sigC & CExtraMask) != UInt(0)) ^ doSubMags)(sigSumSize-1, 0)
```

27

# hardfloat: mulAddSubRecodedFloatN (Expert-written)
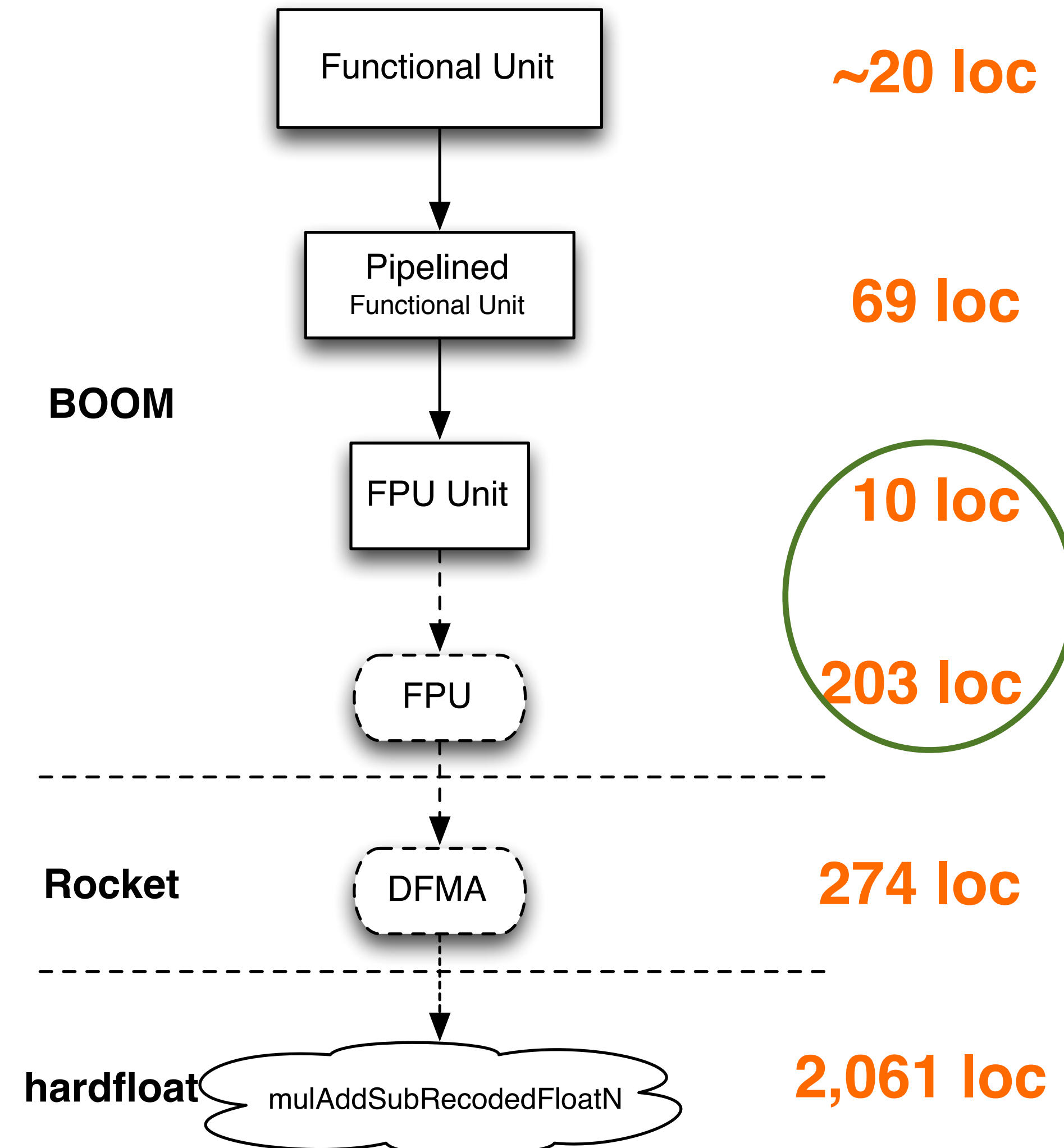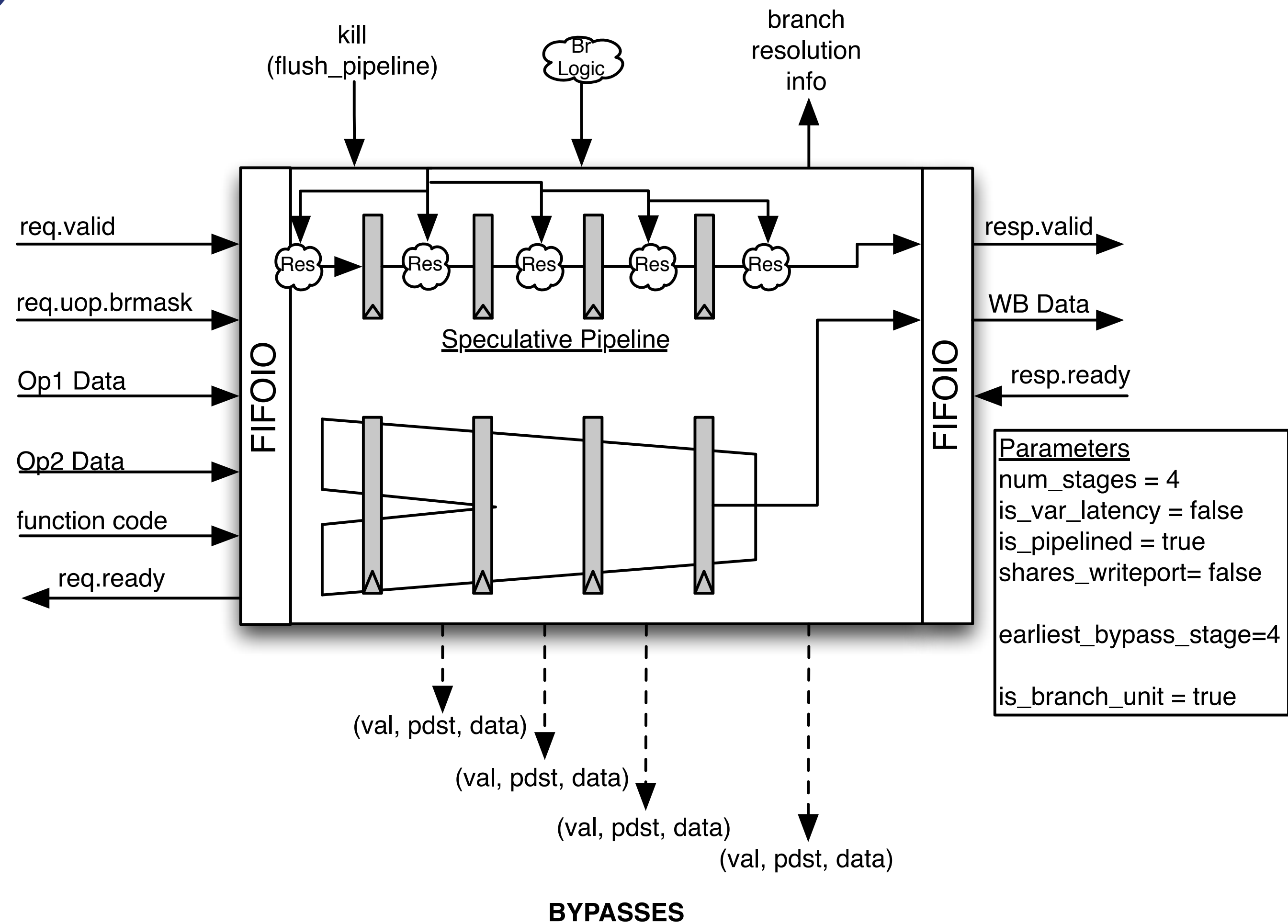
**a little snippet of lines (189-207):**

```scala
val notCDom_signSigSum = sigSum(normSize+1)
val doNegSignSum =
    Mux(isCDominant, doSubMags & ~ isZeroC, notCDom_signSigSum)
val estNormDist =
  Mux(isCDominant, CDom_estNormDist,
  Mux(notCDom_signSigSum, estNormNeg_dist,
  estNormPos_dist))
val cFirstNormAbsSigSum = // ??? odd mux gives the best DC synthesis QoR
  Mux(notCDom_signSigSum,
    Mux(isCDominant, CDom_firstNormAbsSigSum, notCDom_neg_cFirstNormAbsSigSum),
    Mux(isCDominant, CDom_firstNormAbsSigSum, notCDom_pos_firstNormAbsSigSum))
val doIncrSig = ~ isCDominant & ~ notCDom_signSigSum & doSubMags
val estNormDist_5 = estNormDist(logNormSize-3, 0).toUInt
val normTo2ShiftDist = ~ estNormDist_5
val absSigSumExtraMask = Cat(MaskOnes(normTo2ShiftDist, 0, firstNormUnit-1), Bool(true))
val sigX3 =
    Cat(cFirstNormAbsSigSum(sigWidth+firstNormUnit+3,1) >> normTo2ShiftDist,
     Mux(doIncrSig, (~cFirstNormAbsSigSum(firstNormUnit-1,0) & absSigSumExtraMask) === UInt(0),
      (cFirstNormAbsSigSum(firstNormUnit-1,0) & absSigSumExtraMask) != UInt(0)))(sigWidth+4, 0)
```

## and at the top...

```
//*** THIS MODULE HAS NOT BEEN FULLY OPTIMIZED.
//*** DO THIS ANOTHER WAY?
```

# Adding Floating-point



branch resolution info

kill (flush_pipeline)

Br Logic

req.valid

req.uop.brmask

Op1 Data

Op2 Data

function code

req.ready

resp.valid

WB Data

resp.ready

Speculative Pipeline

FIFOIO

FIFOIO

Res  Res  Res  Res  Res

Parameters
num_stages = 4
is_var_latency = false
is_pipelined = true
shares_writeport= false

earliest_bypass_stage=4

is_branch_unit = true

(val, pdst, data)
(val, pdst, data)
(val, pdst, data)
(val, pdst, data)

**BYPASSES**

Functional Unit          ~20 loc

Pipelined
Functional Unit          69 loc

**BOOM**

FPU Unit                 10 loc

FPU                      203 loc

**Rocket**    DFMA        274 loc

**hardfloat**  mulAddSubRecodedFloatN    2,061 loc

- ■ 12 days to add SP, DP floating point
- ■ 1092 lines of code added

**29**

# Adding Floating-point

```scala
class FPUUnit(num_stages: Int) extends PipelinedFunctionalUnit(
    num_stages = num_stages,
    num_bypass_stages = 0,
    earliest_bypass_stage = 0,
    data_width = 65)
    with BOOMCoreParameters
{

    val fpu = Module(new FPU())
    fpu.io.req <> io.req
    fpu.io.req.bits.fcsr_rm := io.fcsr_rm
    io.resp <> fpu.io.resp
    io.resp.bits.fflags.bits.uop := io.resp.bits.uop

}
```

Functional Unit — ~20 loc

Pipelined Functional Unit — 69 loc

BOOM

FPU Unit — 10 loc

FPU — 203 loc

dumb, compute pipe

Rocket — DFMA — 274 loc

hardfloat — mulAddSubRecodedFloatN — 2,061 loc

kill (flush_pipeline)

Br Logic

branch resolution info

req.valid

req.uop.brmask

Op1 Data

Op2 Data

function code

req.ready

FIFOIO

Res  Res  Res  Res  Res

Speculative Pipeline

FIFOIO

resp.valid

WB Data

resp.ready

Parameters
num_stages = 4
is_var_latency = false
is_pipelined = true
shares_writeport= false

earliest_bypass_stage=4

is_branch_unit = true

(val, pdst, data)
(val, pdst, data)
(val, pdst, data)
(val, pdst, data)

**BYPASSES**

# Load/Store Unit

# Building a Register File (the first P&R)

- BOOMv1 -- 7r3w with 110 registers (INT/FP)
- Initial Regfile design was infeasible for layout
- critical paths in issue-select and register read
- Not DRC/LVS clean

# Multi-port Register File for Design Exploration

## Transistor-level



**Advantage**

- Compact area
- Higher performance

**Challenge**

- Long design cycle
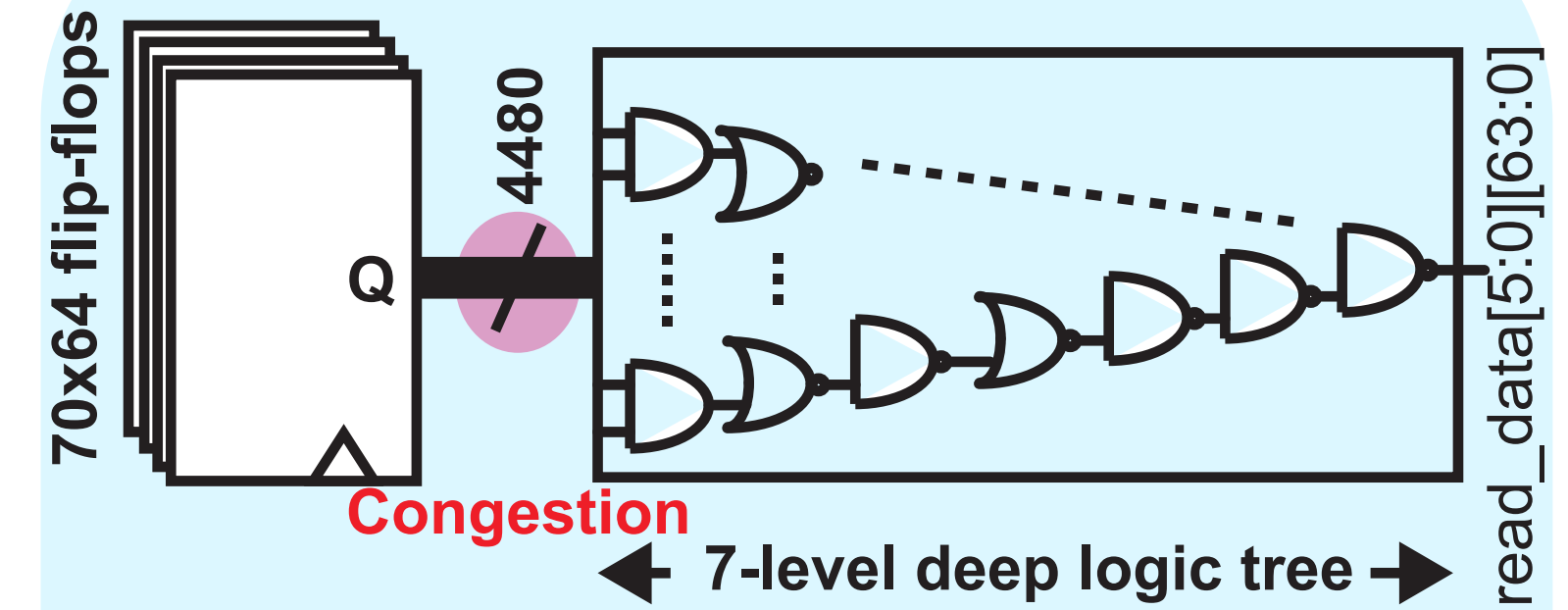- Difficult for architecture design exploration

## Gate-level



**Advantage**

- Rapid design exploration
- Shared read wires solve routing congestion

**Challenge**

- Guided place-and-route for area/performance optimization

## RTL



**Advantage**

- Low design effort
- Rapid design exploration

**Challenge**

- Large area
- Bad performance
- Routing congestion

# Verification

- Directed tests and a randomized torture generator (riscv-torture).
- Verilator/VCS/FPGA simulation at RTL.
- VCS for post-gl/par simulation.
- Speculative OOO pipelines are difficult to get good coverage on.
  - Need tests that build up a lot of speculative state.
  - Need tests that cover OS- and platform-level use-cases.
- Assertions are king.
- Currently moving towards using **co-simulation** against an ISA simulator (using CSRFile's trace port).

# Using the Code

- https://github.com/riscv-boom/riscv-boom
- IntelliJ is awesome
  - https://github.com/riscv-boom/boom-template#vimbash-isnt-a-development-environment-how-do-i-setup-an-intellij-ide

# Code Repo Organization



- rocket-chip, riscv-boom are git submodules of boom-template
- boom-template is just a template for gluing an SoC together
- source code lies in riscv-boom
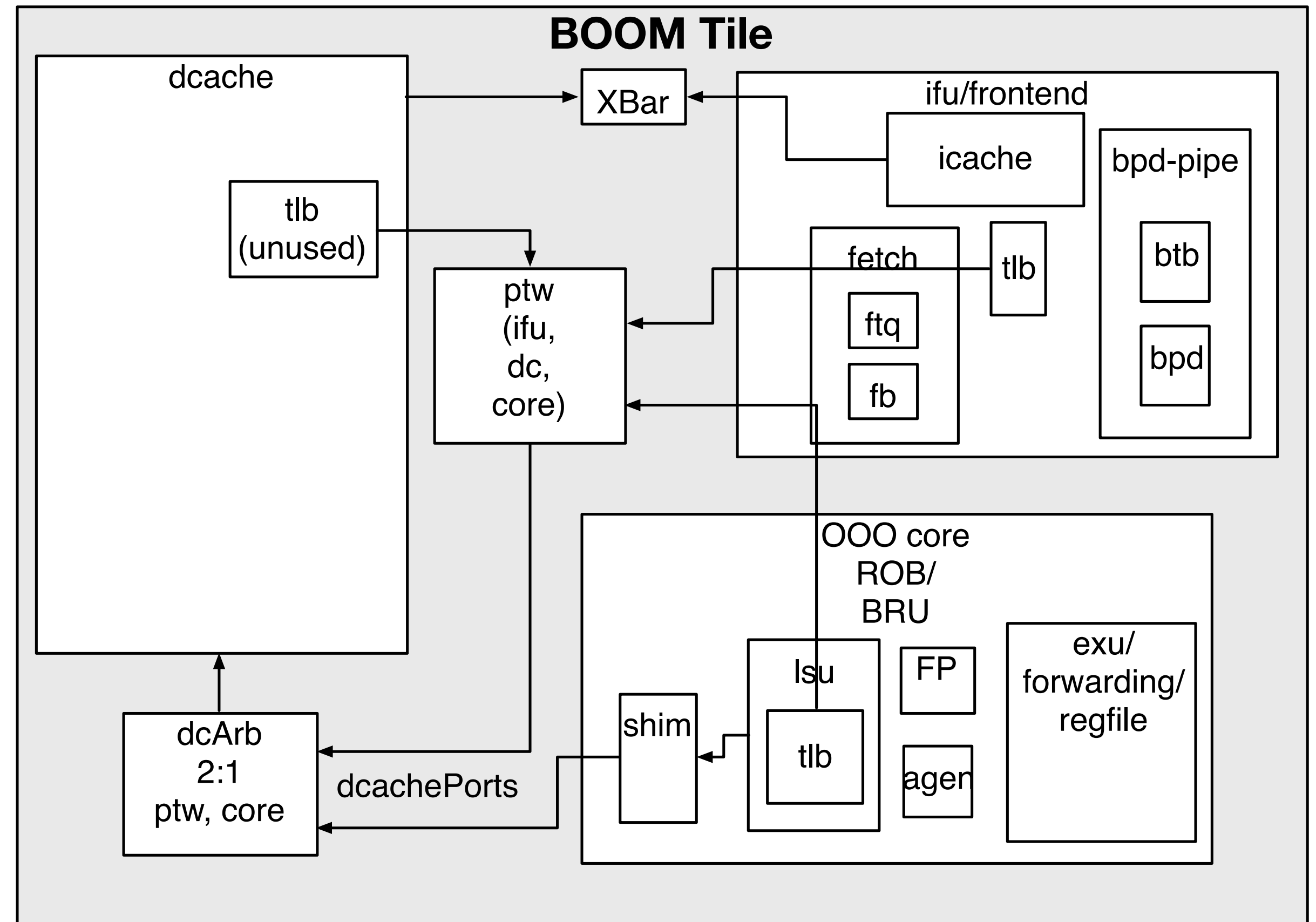- rocket-chip is a library, provides the uncore

# Code Repo Organization



- boom-template is just a template for gluing an SoC together
- fork template for tape-outs

- ■ From rocket-chip:
  - – dcache, icache, tlb, ptw, XBar, dcArb
  - – caches have been hard forked for BOOM-specific changes (e.g., Spectre-related)
- ■ Historical artifact:
  - – OOO core was "BOOM", IFU and DC came from rocket-chip
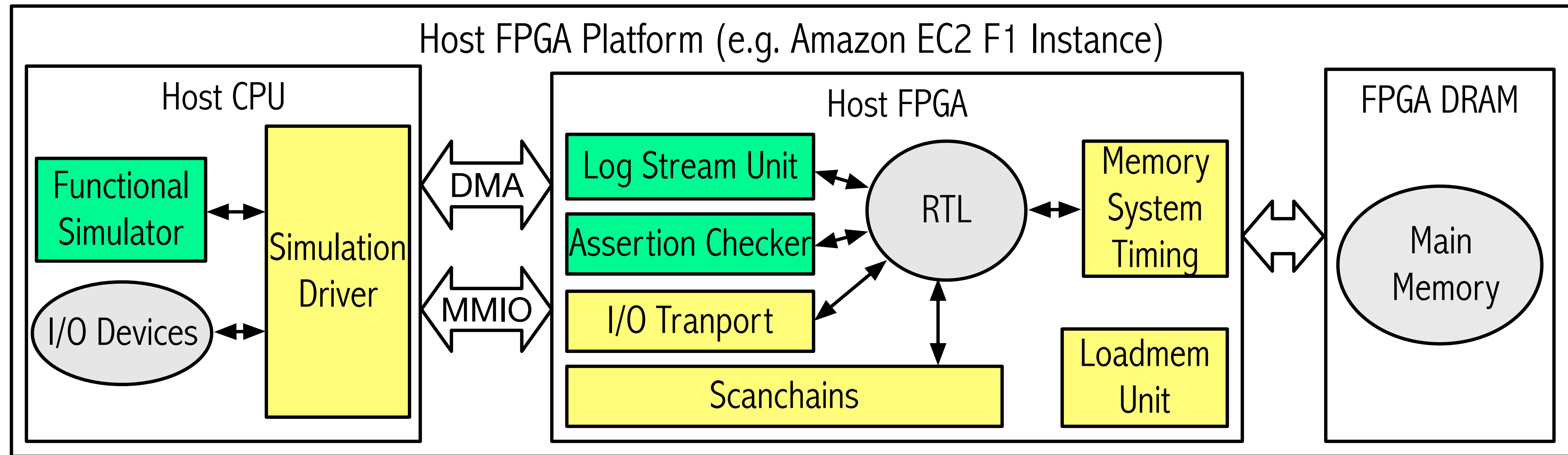
# Questions before the transition?

# Extra Slides

# DESSERT: Debugging RTL Effectively with State Snapshotting for Error Replays across Trillions of Cycles



Host FPGA Platform (e.g. Amazon EC2 F1 Instance)

Host CPU: Functional Simulator, Simulation Driver, I/O Devices

DMA, MMIO

Host FPGA: Log Stream Unit, Assertion Checker, I/O Tranport, Scanchains, RTL, Memory System Timing, Loadmem Unit

FPGA DRAM: Main Memory

⬤ : Target Module    🟨 :Existing Simulation Component    🟩 : Debugging Module

- Co-simulate, find bugs, and get waveforms from Cloud FPGA-based simulation!
- Donggyu Kim, et. al. CARRV 2018
- https://carrv.github.io/2018/papers/CARRV_2018_paper_10.pdf

# Incorrect Jump Target

- 401.bzip2 (assertion error at 500 billion cycles)
  - JAL jumps to wrong target.
  - Due to improper signed arithmetic.
  - 2-3 year old bug.
  - 3 hours of FPGA time.
  - Would require 39 years of Verilator simulation to find.
  - DESSERT found this via a synthesized assertion.

- 445.gobmk (assertion error at 14.9 billion cycles)
  – misspeculated FPtoInt writes back to invalid ROB entry (after being killed)
  – introduced when splitting regfile into separate integer and fp regfiles
- Problem
  – FPtoInt moves share write port with loads
  – FPtoInt gets buffered in a queue; then gets killed
  – queue then later writes back the value anyways when next FPtoInt instruction comes in
- Cause
  – copy/paste error from a non-speculative flow-through queue
- Found
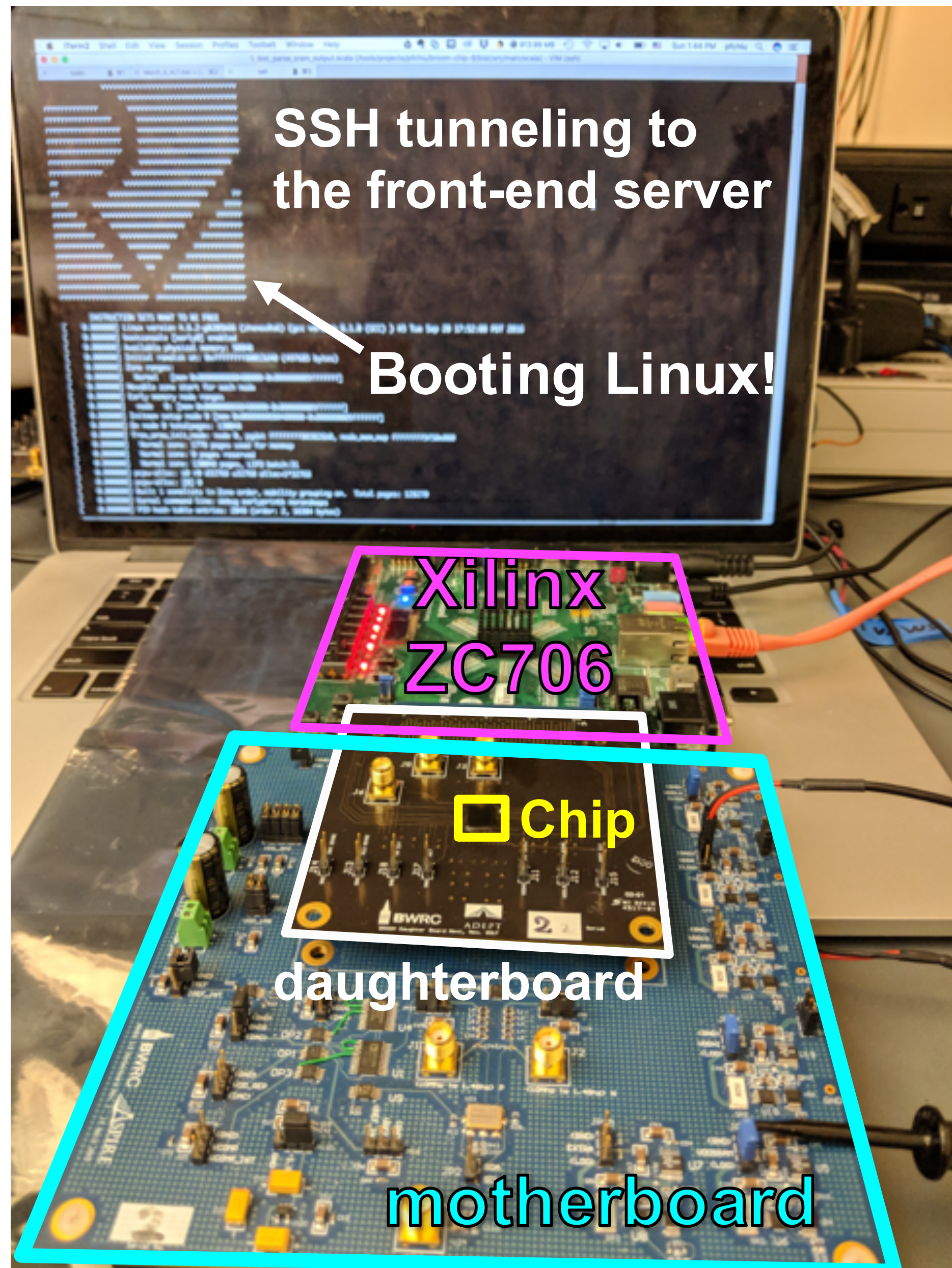  – fpga simulation (as opposed to 431 days of verilator)

# Load-reserve Didn't Page Fault

- LR is a member of the Atomics Extension
- Thought it was an AMO -- so ignored load page fault signal
- Returned garbage from memory
- 4 year old bug
- But... all LRs are followed by a Store-conditional
- The SC
  - takes the page fault
  - fails to get reservation on the first try
  - forces a retry of the LR
  - The LR succeeds and gets the correct data!

# Experimental Setup



- Chip-on-board (COB) package
- Voltage and clock generation on the motherboard
- Cortex A9 on ZC706 works as the front-end server
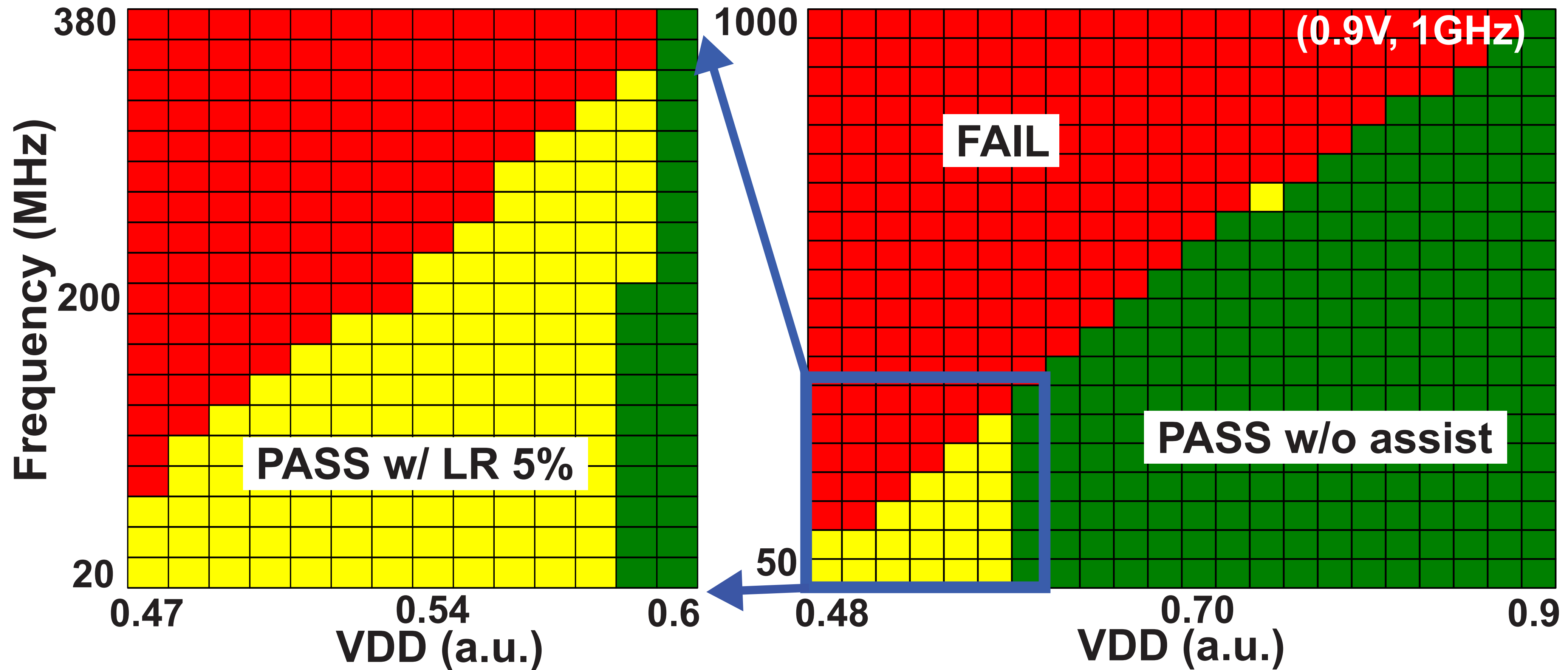- Boot Linux

| Performance | |
|---|---|
| Clock frequency | 1GHz @0.9V |
| | 320MHz @0.6V |
| Coremark/MHz | 3.77 |
| Instruction Per Cycle | 1.11 (@Coremark) |

# Operating voltage and frequency

Benchmark: vvadd



- With LR and 5% loss of L2 cache capacity, Vmin is reduced to 0.47V@70MHz
- 2.3% increase in L2 misses, but only 0.2% degradation in IPC

- Nope, I'm fine actually.
- A TLB permission escalation error
  - allows user to speculatively execute on supervisor data
  - speculative execution leaks information
- In BOOM, TLB permission check fails immediately and squashes load-data bypass
  - no speculative user-level execution using privileged data

# Spectre

- Uh oh
- In BOOM, BTB is shared across threads.
  - allows attacker to force a victim to execute malicious gadget
  - need to flush BTB on context switch, etc.
- What if attacker thread is the victim thread?
  - Uh oh
- What if the attacker invokes a syscall with untrusted input which leaks speculative information?
  - Uh oh
- Great opportunity for security research!
- See my thoughts on this here:
  - https://content.riscv.org/wp-content/uploads/2018/05/13.00-13.15-Celio-Barcelona-Workshop-8-Talk.pdf