# Model based system configuration and tasteful hardware

**Reto Achermann,** Lukas Humbel, Daniel Schwyn, David Cock and Timothy Roscoe

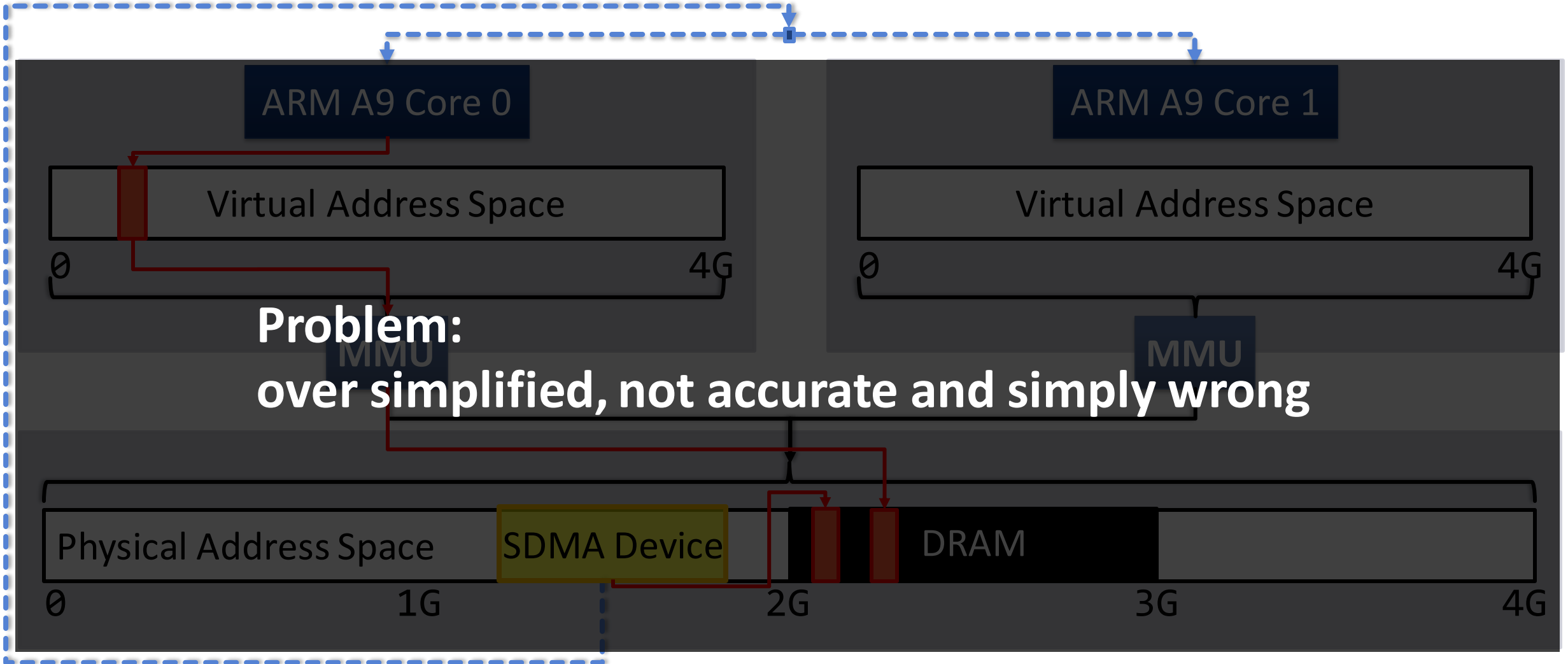Systems Group, Department of Computer Science, ETH Zurich

**Tasteful hardware – or the pursuit of happiness for a systems programmer**

- Operating systems **wrongly** assume a single, uniform address space

- DeviceTrees fail to express important **details** of the hardware configuration

- Programming devices imposes a lot of **constraints** on resource management (memory regions, interrupts, …)

Ideally, we like to accurately **express** hardware configuration and then **generate** a configuration for it.

# Tasteful hardware – or the textbook abstraction



**Problem:** over simplified, not accurate and simply wrong

ARM A9 Core 0

ARM A9 Core 1

Virtual Address Space

0          4G

Virtual Address Space

0          4G

MMU

MMU

Physical Address Space    SDMA Device    DRAM

0          1G          2G          3G          4G

Ti OMAP 4460 SoC

# Reality: The devil is in the details

*Your mobile phone... 5-10 years ago!*

6+ heterogeneous cores

Private and shared memory

5+ Interconnects

Devices attached to different interconnects

Complex interrupt subsystem
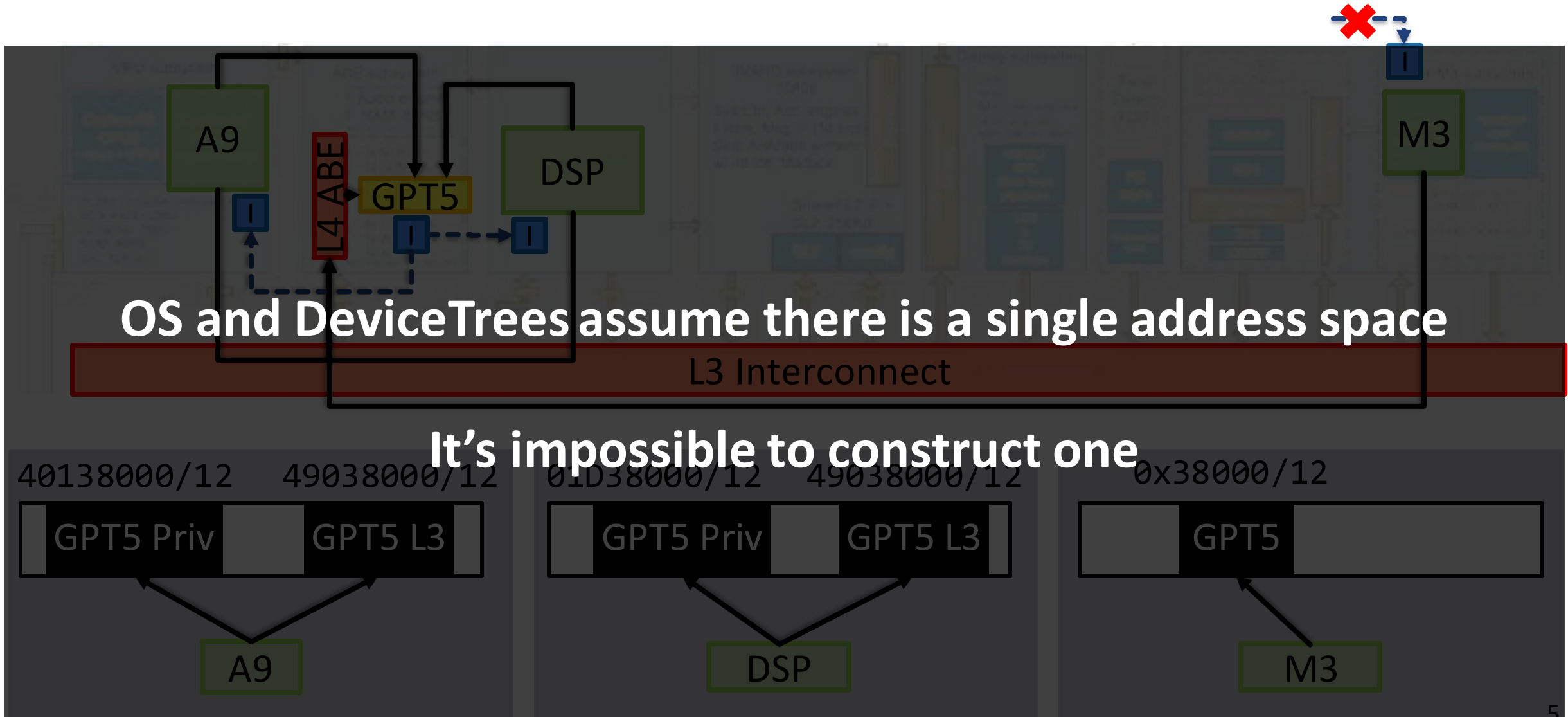
OMAP 4460 SoC, Technical Reference Manual

Next we zoom in here

Takeaway: There are many **details**

# There is NO uniform view of the system



A9

L4 ABE

GPT5

DSP

M3

**OS and DeviceTrees assume there is a single address space**

L3 Interconnect

**It's impossible to construct one**

40138000/12    49038000/12    01D38000/12    49038000/12    0x38000/12

GPT5 Priv    GPT5 L3    GPT5 Priv    GPT5 L3    GPT5

A9    DSP    M3

# We need an accurate representation of hardware

- Experience from the Barrelfish operating system: dealing with this complexity every day.
  - *PCI programming, devices and device drivers,*
  - *Heterogeneity, new architectures, interconnects and platforms*
  - *Resource management: memory allocation, interrupt vectors, …*

  Barrelfish.org

- We don't want to make the same **wrong/imprecise** assumptions

- Provide a sound **hardware description** to facilitate system software **verification**

# System programmer's wishes

- Write **correct** systems code, and do this **fast**
  *Memory management, device drivers, support for new platforms and architectures*

- Make **accurate statements** about hardware platforms
  *At which address can a core/device access memory, what are the access constraints*

- **Generate** system configuration and synthesize configuration algorithms
  *Page tables, DMA bus addresses*

# Outline

1. **Formal model** to accurately represent memory and interrupts subsystems

2. Generate **configuration** based on the model specification

3. Model **refinement** example using MIPS R4600 TLB

## 1. A formal model for memory accesses and interrupts
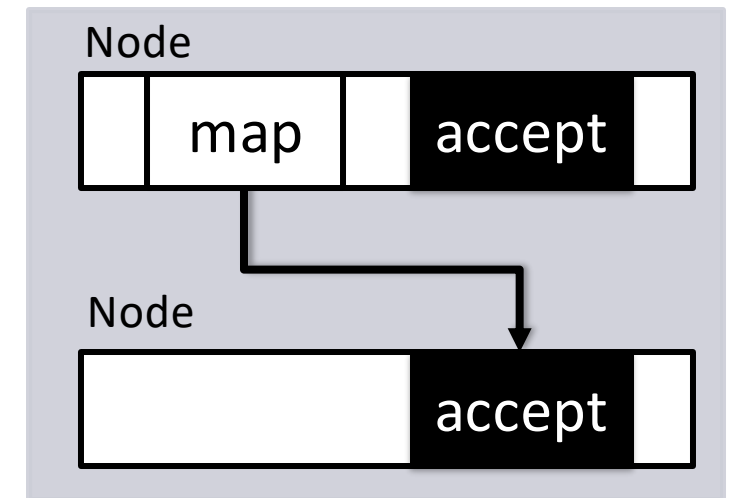
# Design goals of the formal model

- Close to hardware: capture as much **details** as possible

- Represent the **interactions** between the components
  *signal path of a memory access / interrupt*

- Enable different **views/observers** of the system
  *e.g. cores / devices see the system differently*

# Interaction of hardware components on the OMAP4460



Observation: hardware components have well defined **input and output ports.**

➔ Directed graph

# Modelling memory accesses and interrupts as a decoding net

- The model is a **decoding net,** a directed graph

- Hardware components ⇔ **nodes**

  Address spaces, translation units, devices,
  interrupt controllers, interconnects, …

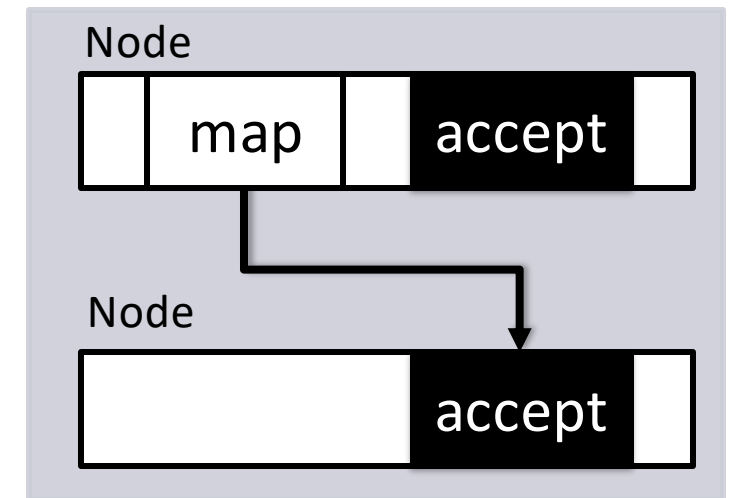- Formally: assignment of identifiers to nodes

$$net: \mathbb{N} \rightarrow node$$

# Modelling memory accesses and interrupts as a decoding net

- Nodes have **two properties**:
  - **Accept**: set of addresses this node responds to

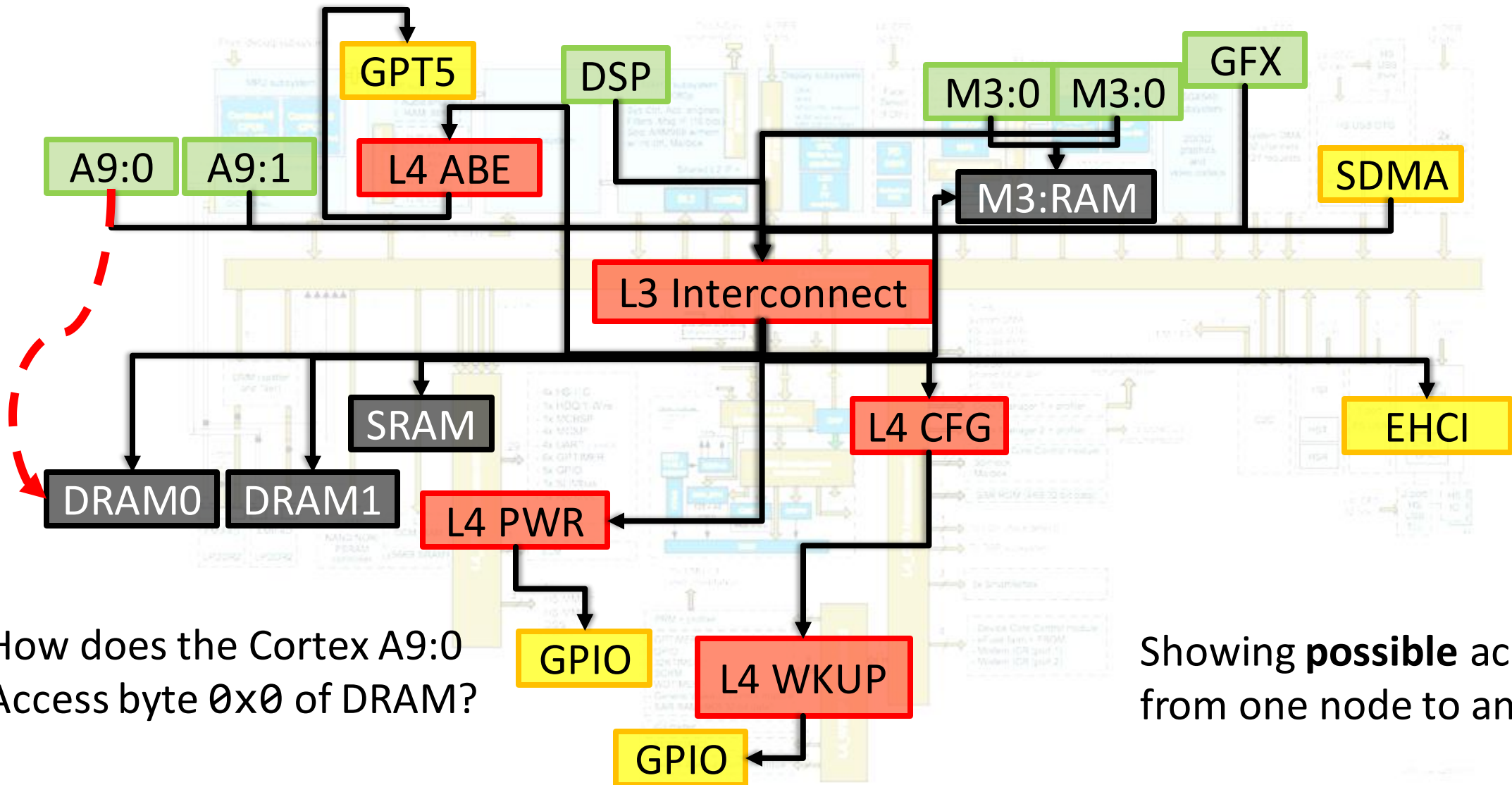  - **Translate**: map input address at node to a name

node:
$$accept :: node \rightarrow \{\mathbb{N}\}$$
$$translate :: node \rightarrow \mathbb{N} \rightarrow \{name\}$$



- Name is a **qualified address** at a node: $name = (\mathbb{N}, \mathbb{N})$

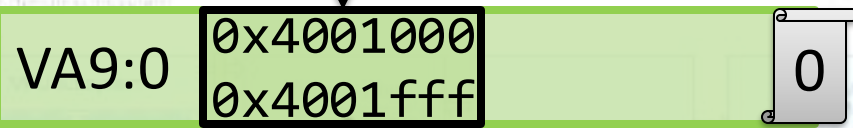# A possible decoding net for the OMAP 4460
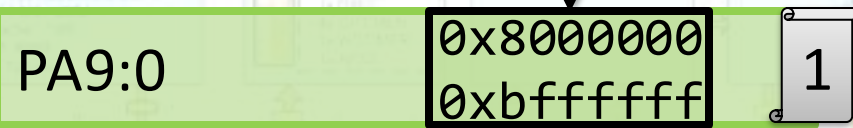


How does the Cortex A9:0
Access byte 0x0 of DRAM?

Showing **possible** accesses
from one node to another

# Modelling the access to byte `0x0` of DRAM from an A9 core

Resolve: `(0, 0x4001000)`

| VA9:0 | `0x4001000` `0x4001fff` | 0 |

`(1, 0x8000000)`

| PA9:0 | `0x8000000` `0xbffffff` | 1 |

`(2, 0x8000000)`

| L3 Interconnect | `0x8000000` `0xbffffff` | 2 |

`(3, 0x0)`

| `0x0000000` `0x4000000` | DRAM | 3 |

Accept `0x0`

Each node has a label

Resolve a **name** (node, address)
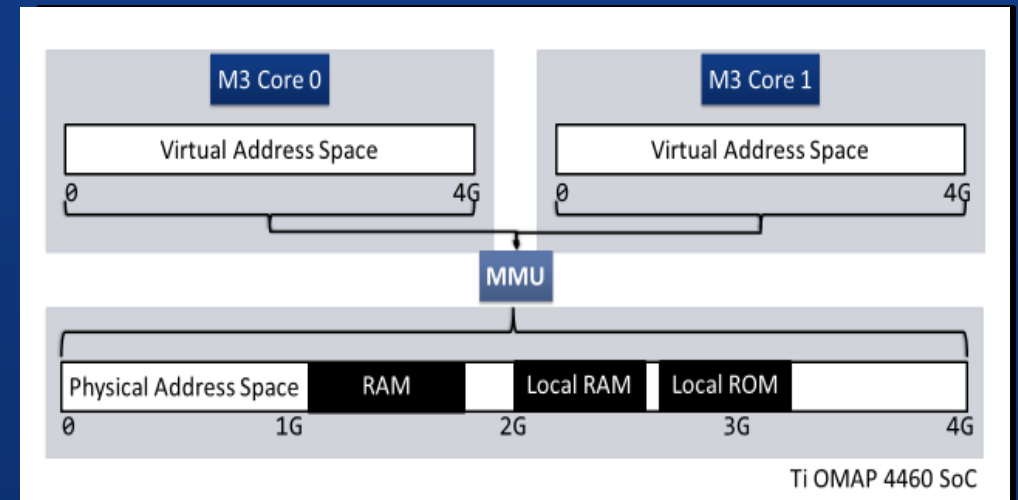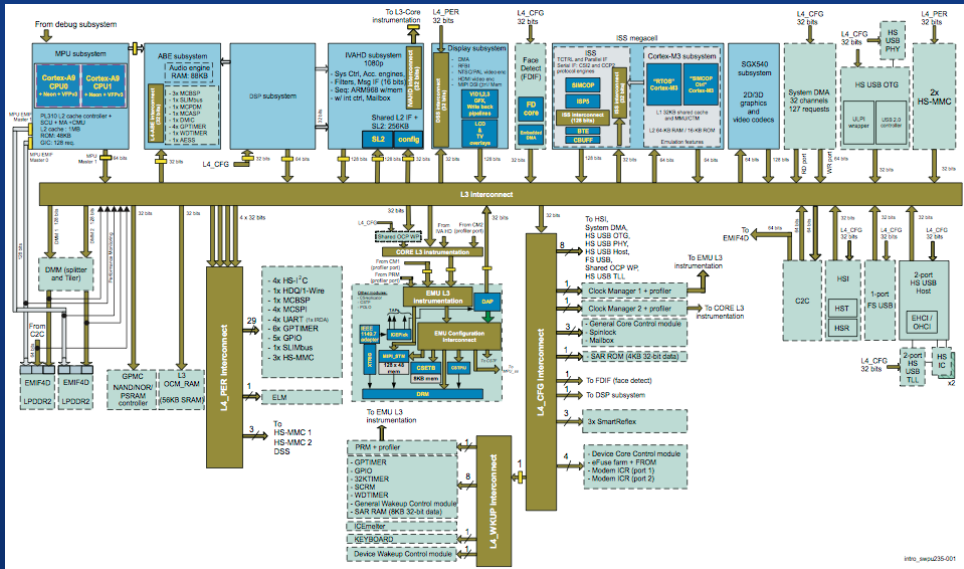
**Model of one particular, static configuration state**

# View Equivalence

Resolve: (0, 0x4001000)                    Resolve: (0', 0x4001000)



**For <u>ONE</u> observer the flattened representation is <u>equivalent</u> to the textbook abstraction**
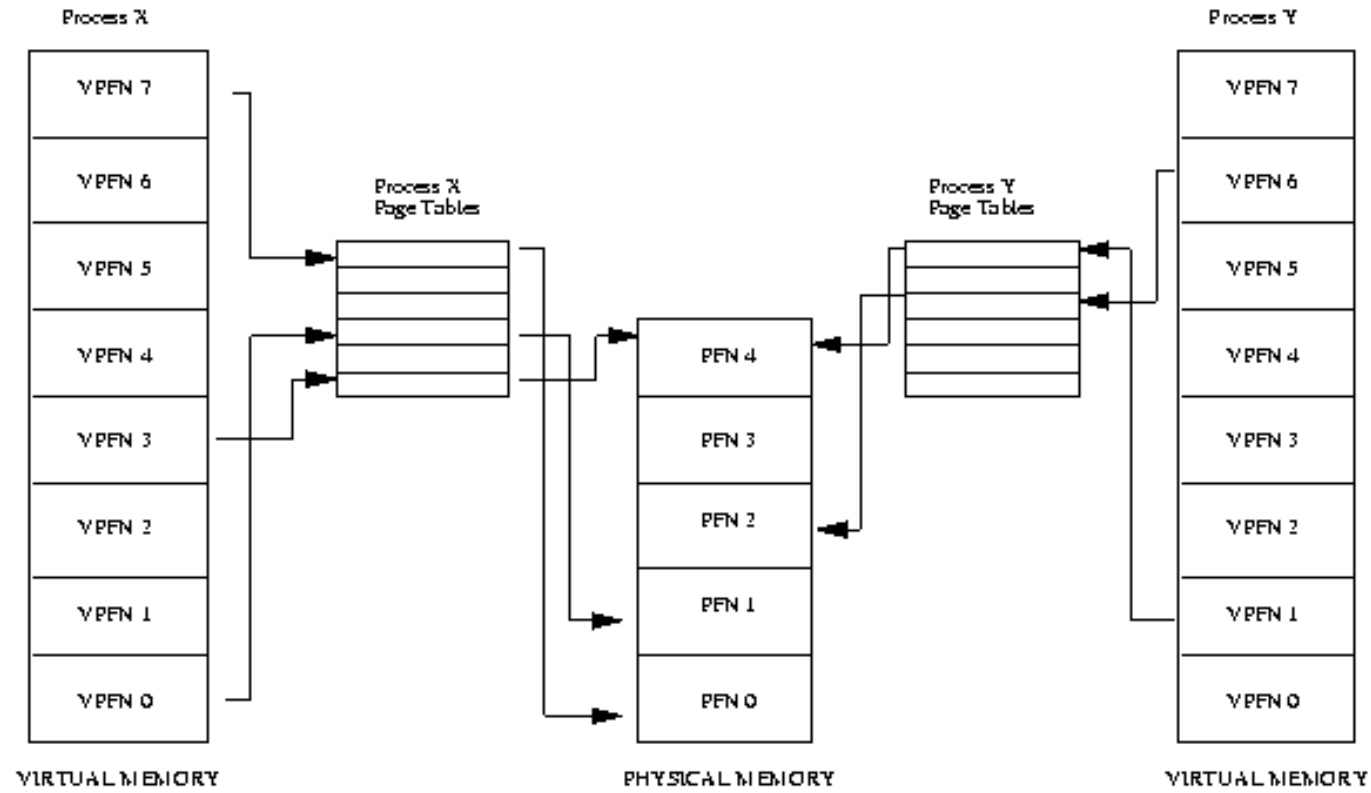
0x4000000    DRAM

Accept 0x0

# Summary

- Express the hardware configuration as a **decoding net model**

- View equivalence preserving transformations for a **fixed observer**

- Formally proven using Isabelle/HOL

- Memory access – Interrupt duality:
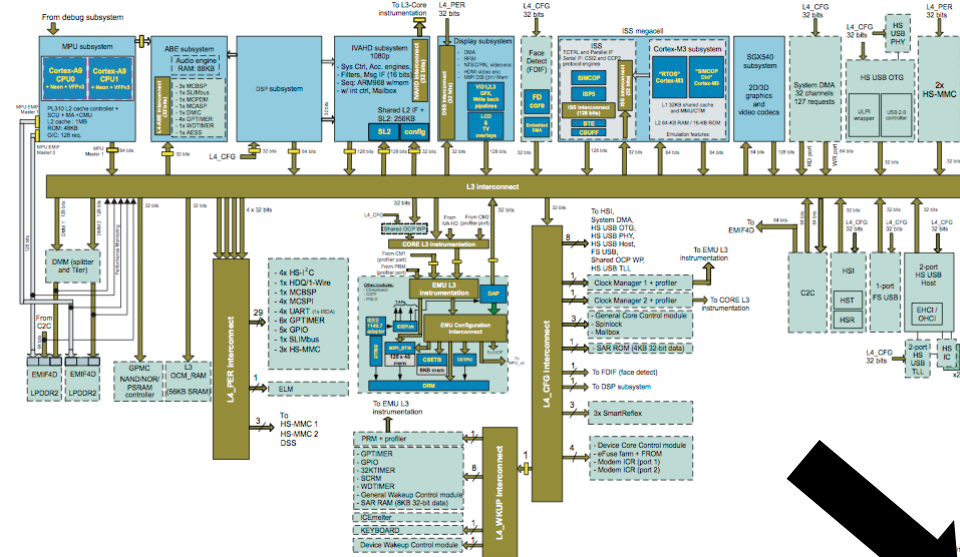  The model is also applicable to the Interrupt subsystem

# 2. Sockeye: Generate Hardware Configurations

Daniel Schwyn, Master's Thesis

## Application of the model in the context of SoC and hardware discovery

- At boot, the kernel needs to know about memory, devices, cores etc.
  ➔ **Devices trees** assumes a **uniform view**

- Our formal model accurately **describes hardware**
  ➔ generate hardware configuration based on the model (*e.g. page tables*)
  ➔ use description for runtime decisions (*e.g. memory allocation*)

- **Sockeye: a domain specific language to describe the hardware configuration**

# Sockeye: a DSL for describing hardware



Models for x86 desktop & servers, Xeon Phi, InfiniBand Cluster, Intel's Single Chip Cloud Computer

Sockeye: DSL + prolog queries

$V_{A9:0}$ **is map** $[20000_3/12$ **to** $P_{A9:0}$ **at** $80000_3]$ $\qquad$ $V_{A9:1}$ **is map** $[20000_3/12$ **to** $P_{A9:1}$ **at** $80000_3]$

$P_{A9:0}, P_{A9:1}$ **are map** $[40138_3/12$ **to** $GPT$ **at** $0]$ **over** $L3$ $\qquad$ $V_{DSP}$ **is over** $P_{DSP}$

$P_{DSP}$ **is map** $[1d3e_3/12$ **to** $GPT$ **at** $0]$ **over** $L3$ $\qquad$ $L2_{M3}$ **is map** $[0_{30}$ **to** $L3$ **at** $80000_3]$

$V_{M3}, V_{M3}$ **are over** $L1_{M3}$ $\qquad$ $L1_{M3}$ **is map** $[0_{28}$ **to** $MIF]$

$RAM_{M3}$ **is accept** $[55020_3/16]$ $\qquad$ $L4$ **is map** $[49038_3/12$ **to** $GPT$ **at** $0]$

$ROM_{M3}$ **is accept** $[55000_3/14]$ $\qquad$ $GPT$ **is accept** $[0/12]$

$MIF$ **is map** $[0 - 5\text{ffffff}$ **to** $L2_{M3}, 55000_3/14$ **to** $RAM_{M3}, 55020_3/16$ **to** $ROM_{M3}]$

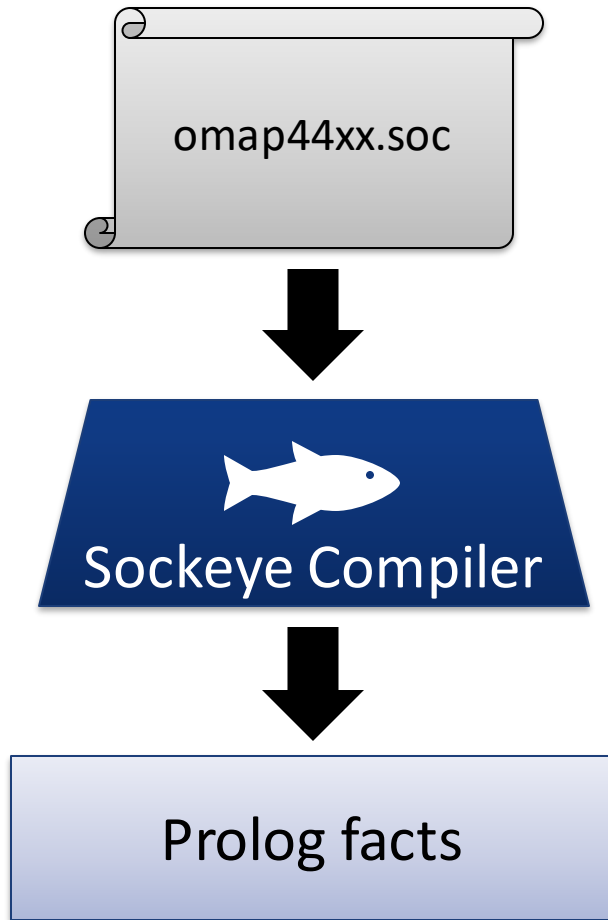$L3$ **is map** $[49000_3/24$ **to** $L4$ **at** $40100_3, 55000_3/12$ **to** $MIF]$ **accept** $[80000_3/30]$

Sockeye Code

# Background: Barrelfish's System Knowledge Base

- Observation: resource allocation and hardware configuration can be expressed as a constrained satisfaction problem.
  *e.g. PCI programming, DMA, …*

- Barrelfish's approach: **SKB = Prolog engine + constraint solver**

- Information about the system is stored as **Prolog facts** in the SKB

- **Prolog queries** express allocation policies and configuration constraints

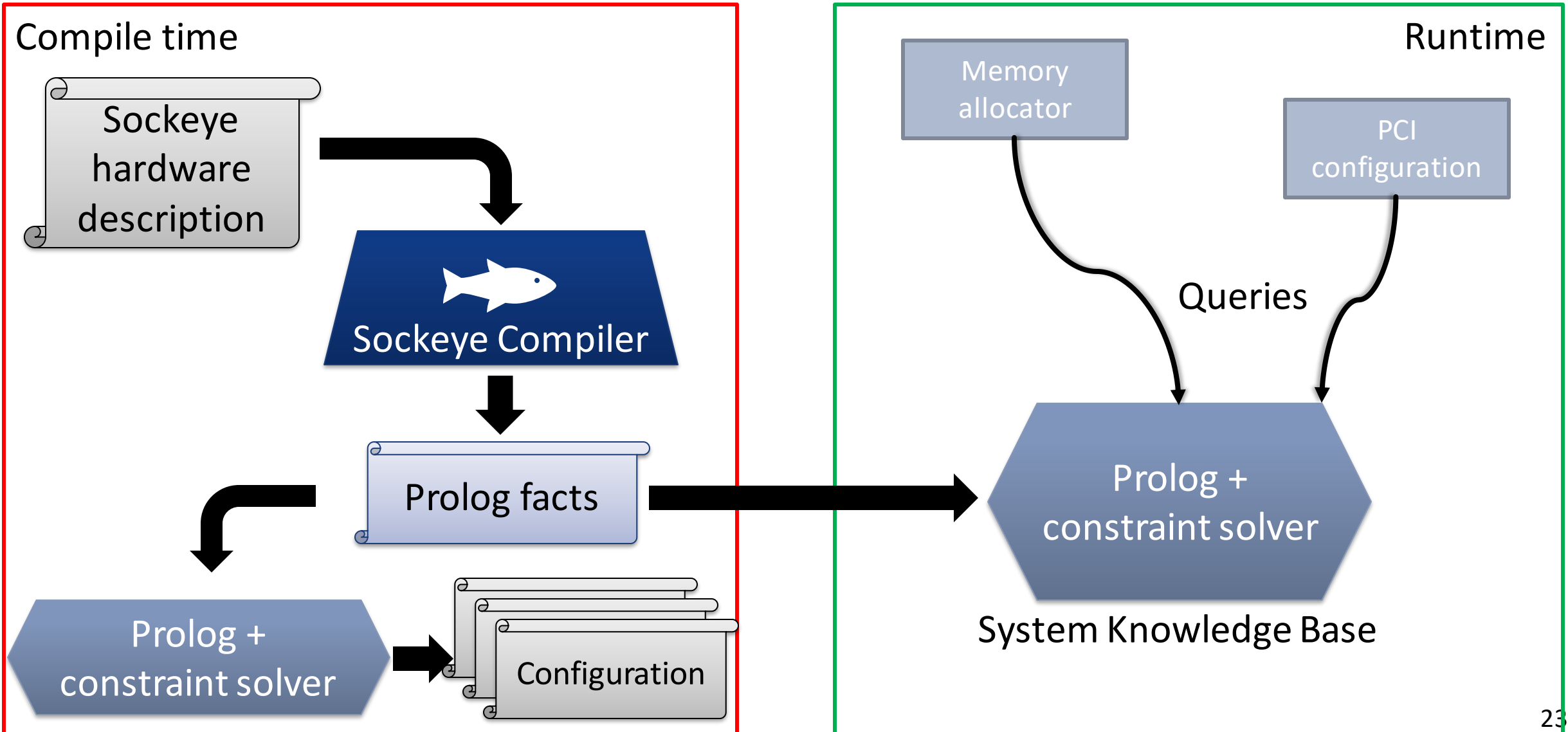# Sockeye Compiler: Generate the decoding net model

omap44xx.soc

$V_{A9:0}$ **is map** $[20000_3/12$ **to** $P_{A9:0}$ **at** $80000_3]$     $V_{A9:1}$ **is map** $[20000_3/12$ **to** $P_{A9:1}$ **at** $80000_3]$

$P_{A9:0}, P_{A9:1}$ **are map** $[40138_3/12$ **to** $GPT$ **at** $0]$ **over** $L3$     $V_{DSP}$ **is over** $P_{DSP}$

$P_{DSP}$ **is map** $[1d3e_3/12$ **to** $GPT$ **at** $0]$ **over** $L3$     $L2_{M3}$ **is map** $[0_{30}$ **to** $L3$ **at** $80000_3]$

$V_{M3}, V_{M3}$ **are over** $L1_{M3}$     $L1_{M3}$ **is map** $[0_{28}$ **to** $MIF]$

$RAM_{M3}$ **is accept** $[55020_3/16]$     $L4$ **is map** $[49038_3/12$ **to** $GPT$ **at** $0]$

$ROM_{M3}$ **is accept** $[55000_3/14]$     $GPT$ **is accept** $[0/12]$

$MIF$ **is map** $[0 - 5fffffff$ **to** $L2_{M3}, 55000_3/14$ **to** $RAM_{M3}, 55020_3/16$ **to** $ROM_{M3}]$

$L3$ **is map** $[49000_3/24$ **to** $L4$ **at** $40100_3, 55000_3/12$ **to** $MIF]$ **accept** $[80000_3/30]$

Sockeye Compiler

Sockeye compiler converts the description into the **decoding net model** expressed as Prolog facts

Prolog facts

```
net('SRAM',node(memory,[block(16'0,16'3fffffff)],[],'@none')).
net('BOOT_ROM',node(memory,[block(16'0,16'bfff)],[],'@none')).
net('L3_OCM_RAM',node(memory,[block(16'0,16'dfff)],[],'@none')).
net('SDRAM',node(memory,[block(16'0,16'3fffffff)],[],'@none')).
```

# Sockeye Runtime Environment

Compile time

Runtime

Sockeye hardware description

Sockeye Compiler

Prolog facts

Prolog + constraint solver

Configuration

Memory allocator

PCI configuration

Queries

Prolog + constraint solver

System Knowledge Base

# Application: Generation of kernel page tables

```
kernelPageTable = pageTableWith {
    cpu = "CORTEXA9",
    mainMemory = "SDRAM",
    devices = [ "UART3", "SCU"]
},
```

Prolog facts
generated by Sockeye

Kernel page
table generator

```
union arm_l1_entry l1_table[ARM_L1_MAX_ENTRIES]
    __attribute__((aligned(ARM_L1_ALIGN), section(".boot.tables"))) =
{
        [L1_TABLE_INDEX(0x7FE00000)] = L1_DEVICE_ENTRY(0x7FE00000),
        [L1_TABLE_INDEX(0x7FF00000)] = L1_DEVICE_ENTRY(0x7FF00000), …
        [L1_TABLE_INDEX(0x80000000)] = L1_MEMORY_ENTRY(0x80000000),
        [L1_TABLE_INDEX(0x80100000)] = L1_MEMORY_ENTRY(0x80100000), …
}
```

# Application: Memory allocation

**DMA** bus address allocation is required to satisfy various **constraints**

A good fit for the SKB:

- The decoding net model can resolve all **reachable memory regions** from the DMA device

- Express the constraints as a prolog query to get **usable memory regions**

# Summary: Configuration generation

- Describe hardware configuration in the Sockeye DSL

- Obtain a Prolog representation of the decoding net model

- Generate hardware configuration

- Using Barrelfish's SKB to perform DMA bus address allocation

# 3. Case Study with the MIPS R4600 TLB

# MIPS R4600 TLB Case study

- Software loaded TLB: explicit management of the TLB entries
  → we don't need to worry about hardware page table walker

- We want to express the MIPS R4600 TLB in the decoding net model

- **Main Question:** Can we refine the abstract decoding net node to a **dynamic** TLB node?

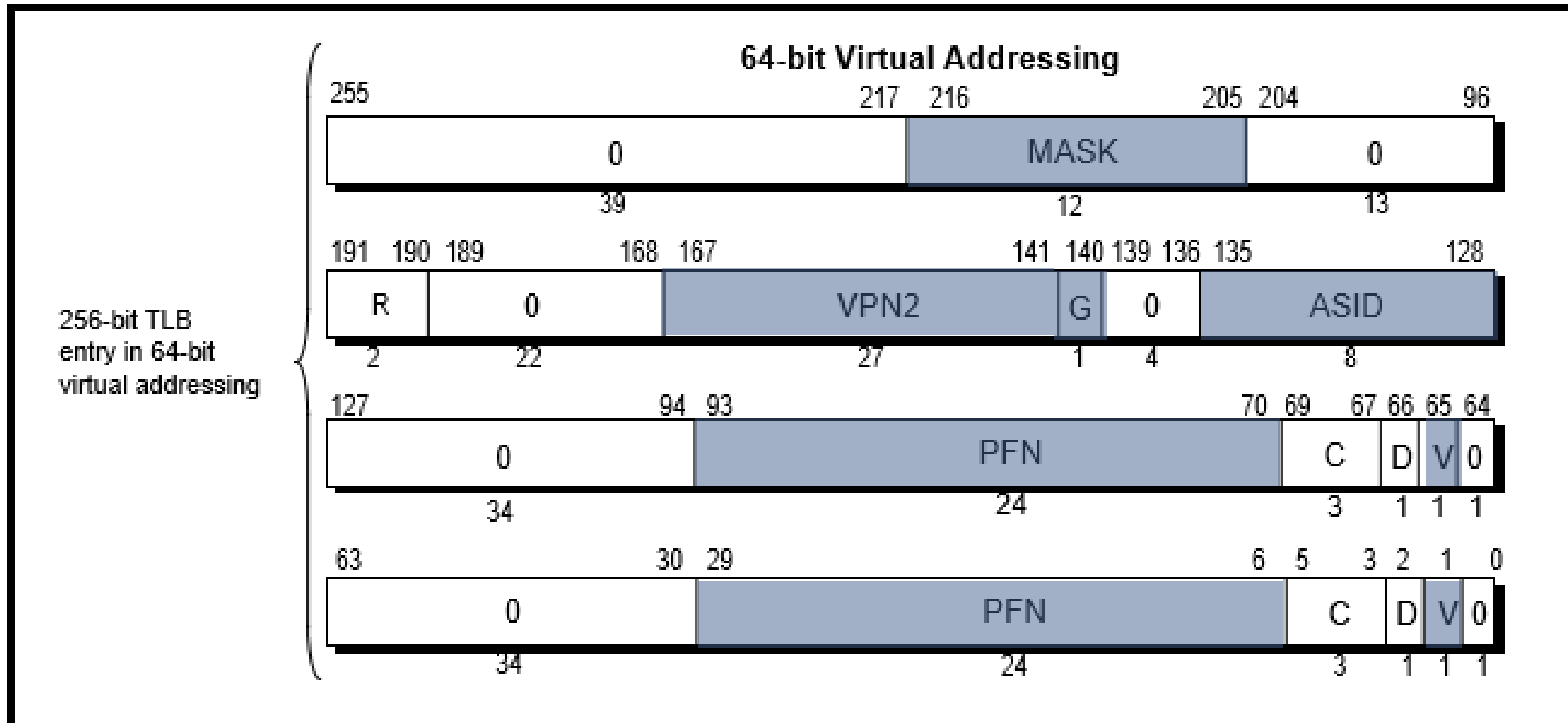# MIPS R4600 TLB entry format at a glance



Figure 4.8  Format of a TLB Entry

# Modelling the MIPS R4600 TLB in Isabelle/HOL

```
(* ------------------------------------------------------------------ *)
subsection "The MIPS TLB"
(* ------------------------------------------------------------------ *)

text "The MIPS TLB has a fixed number of entries ('capacity') where some of them
      are wired and not replaced by the random replacement."

record MIPSTLB =
  wired      :: nat
  entries    :: "nat ⇒ TLBENTRY"

text "The MIPS TLB has in total 48 entries, of which are at maximum 32 wired."

definition TLBCapacity :: nat
  where "TLBCapacity = 48"
```
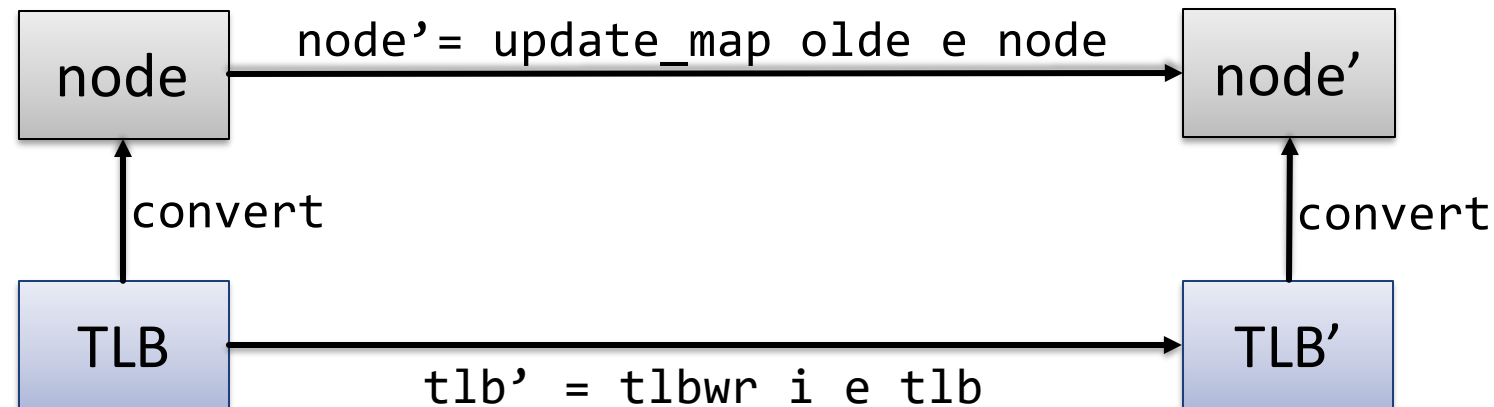
- Formal specification of the four TLB operations with pre and post conditions

- Proofs on operations *e.g. commutativity of updates, …*

# Refinement of the decoding net

- Defined an lifting function from MIPS TLB to decoding net node.

- Formal proof in Isabelle/HOL that an entry update of a well-formed MIPS TLB is reflected at the map function of the decoding net node
  *e.g. write entry e to index i:*

```
               node'= update_map olde e node
   node   ─────────────────────────────────►   node'
    ▲                                            ▲
    │ convert                                    │ convert
    │                                            │
   TLB  ───────────────────────────────────►   TLB'
               tlb' = tlbwr i e tlb
```

# Invariant for the MIPS R4600 TLB

*"The R4600/R4700 does not provide any detection or shutdown mechanism for multiple matches in the TLB. There is no damage possible from this condition. **The result is undefined for this condition. Software is expected never to allow this to occur"***

*-- MIPS R4600 Manual*

**An address must never match two entries at the same time.**

We say a TLB is valid, if all it's entries are well formed and there is at most one entry that matches a particular VPN

**The Invariant cannot be satisfied!**

*"The contents of all registers in the CPU are **undefined** when this [reset] exception occurs, except for the following register fields: [...]"*

*-- MIPS R4600 Manual*

- Therefore: at reset, the TLB entries are undefined and considered random.
- TLB is enabled at boot time, cannot be enabled

- Theorem $\exists t. (TLB\_in\_reset\ t) \land \neg(TLBValid\ t)$
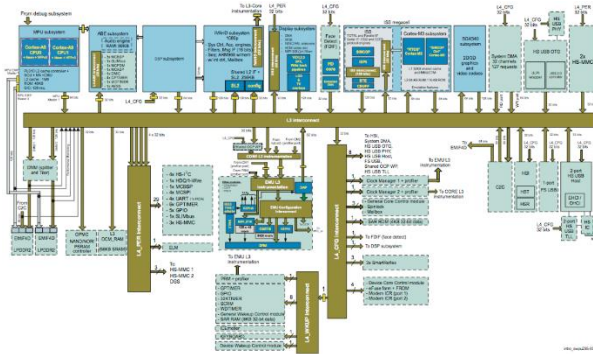  A proof that there exists a TLB in reset which is invalid at the same time

33

# Summary: MIPS TLB Case study

- Formal model of the MIPS R4600 TLB in Isabelle/HOL

- Refinement: from the decoding net node to detailed real hardware

- Invariant cannot be satisfied:
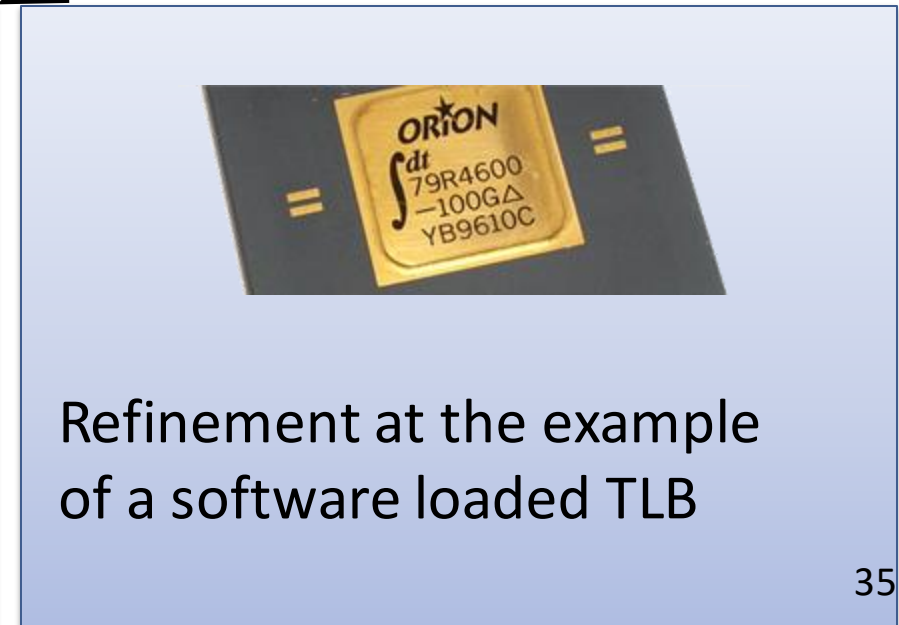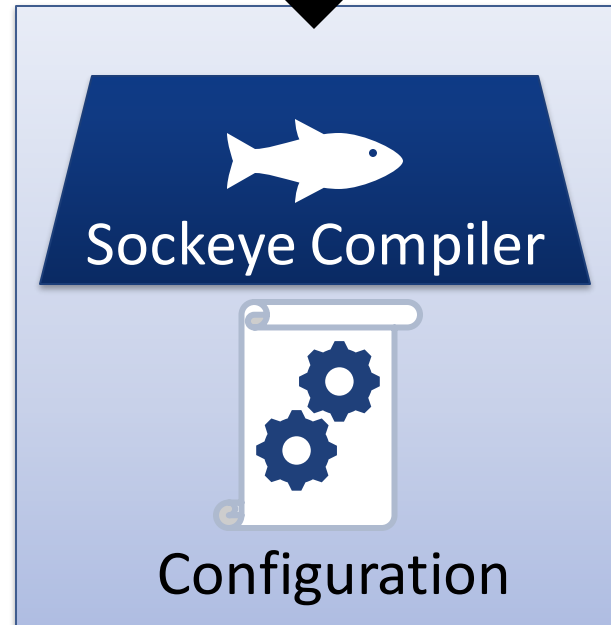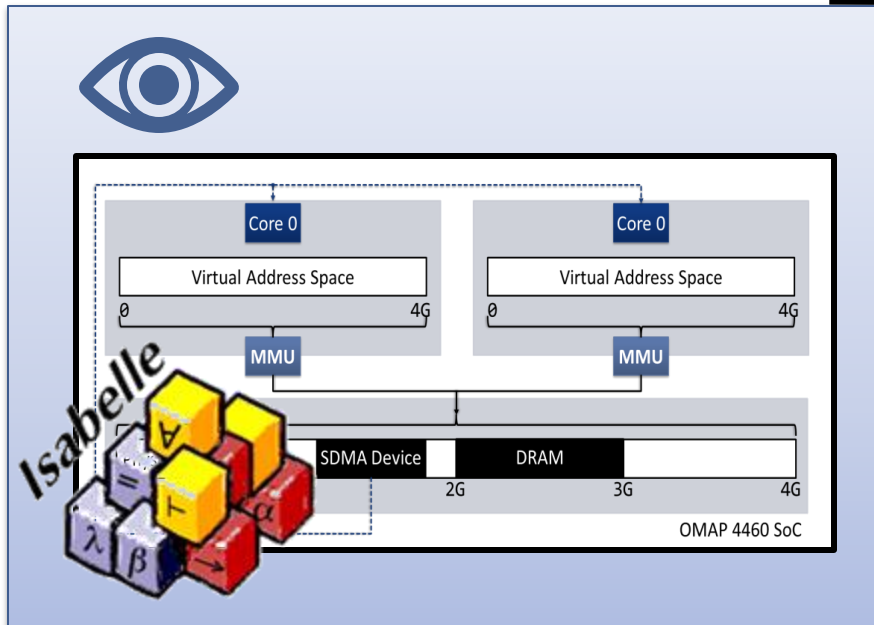  At reset, the TLB is in an undefined (random) state and may contain two matching entries
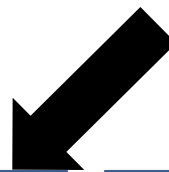
github.com/BarrelfishOS (soon)

# Summary

$$V_{A9:0} \text{ is map } [20000_3/12 \text{ to } P_{A9:0} \text{ at } 80000_3] \quad V_{A9:1} \text{ is map } [20000_3/12 \text{ to } P_{A9:1} \text{ at } 80000_3]$$
$$P_{A9:0}, P_{A9:1} \text{ are map } [40138_3/12 \text{ to } GPT \text{ at } 0] \text{ over } L3 \quad V_{DSP} \text{ is over } P_{DSP}$$
$$P_{DSP} \text{ is map } [1d3e_3/12 \text{ to } GPT \text{ at } 0] \text{ over } L3 \quad L2_{M3} \text{ is map } [0_{30} \text{ to } L3 \text{ at } 80000_3]$$
$$V_{M3}, V_{M3} \text{ are over } L1_{M3} \quad L1_{M3} \text{ is map } [0_{28} \text{ to } MIF]$$
$$RAM_{M3} \text{ is accept } [55020_3/16] \quad L4 \text{ is map } [49038_3/12 \text{ to } GPT \text{ at } 0]$$
$$ROM_{M3} \text{ is accept } [55000_3/14] \quad GPT \text{ is accept } [0/12]$$
$$MIF \text{ is map } [0 - 5fffffff \text{ to } L2_{M3}, 55000_3/14 \text{ to } RAM_{M3}, 55020_3/16 \text{ to } ROM_{M3}]$$
$$L3 \text{ is map } [49000_3/24 \text{ to } L4 \text{ at } 40100_3, 55000_3/12 \text{ to } MIF] \text{ accept } [80000_3/30]$$



Configuration

Refinement at the example of a software loaded TLB

35

# Future work: towards verified configuration



System (re-)configuration



Distinction of Read/Write accesses

- **Synthesize** configuration procedures with transitions between configurations without violation of **invariants**

- **Modularization** of Sockeye

- Distinguish read and write accesses as they do not have the same semantics, read-only