

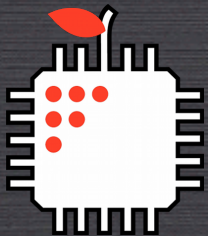
# OPERATING SYSTEMS FOR MANY-CORE

THE PATH FROM TIME SHARING TO SPACE SHARING

RAPHAEL 'KENA' POSS

[R.C.POSS@UVA.NL](mailto:R.C.POSS@UVA.NL)

[WWW.APPLE-CORE.INFO](http://WWW.APPLE-CORE.INFO)



UNIVERSITY OF AMSTERDAM

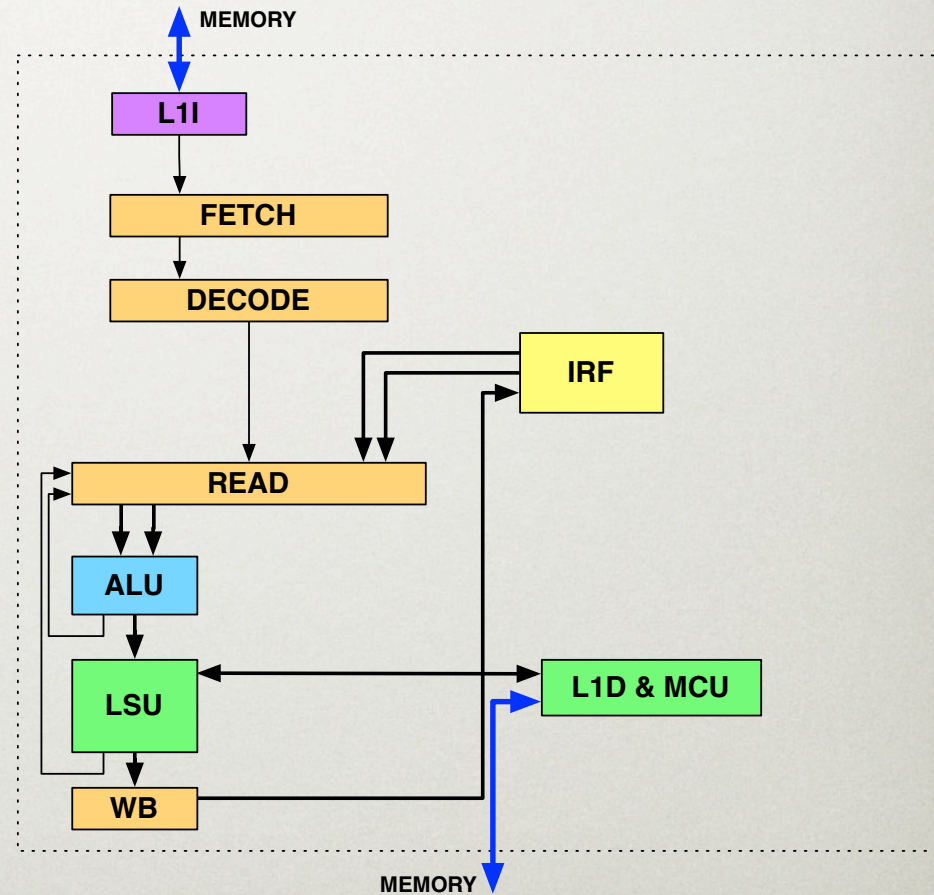




# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

Based on a “simple”  
RISC pipeline...

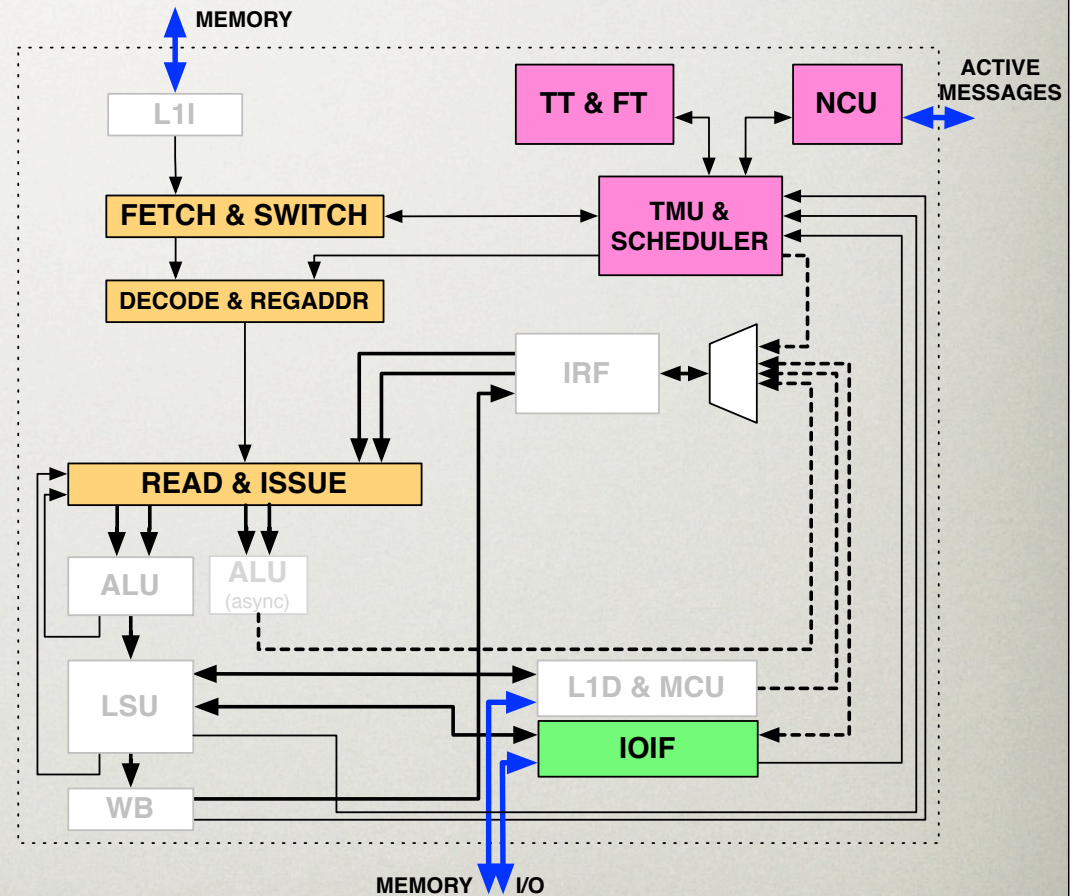




# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

A couple of extensions...



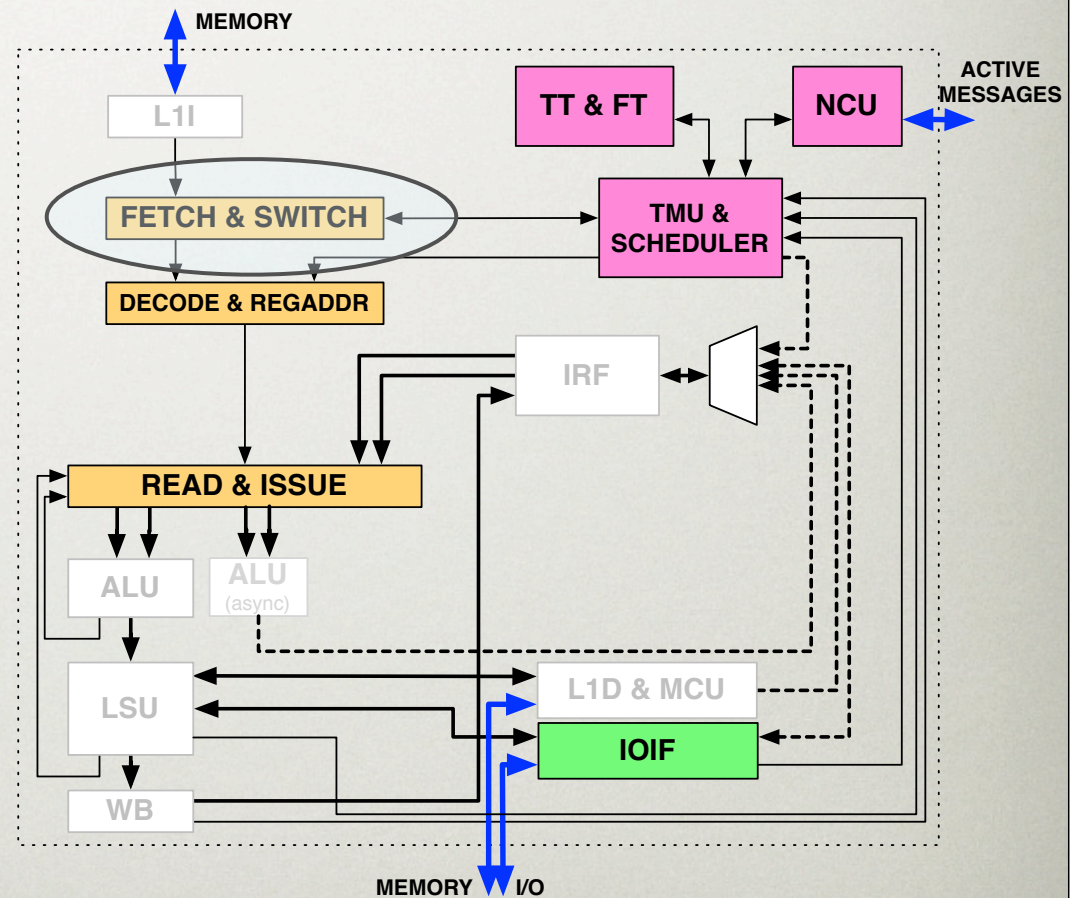


# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

Fetch from FIFO  
of active threads

Switch based on  
instruction annotations



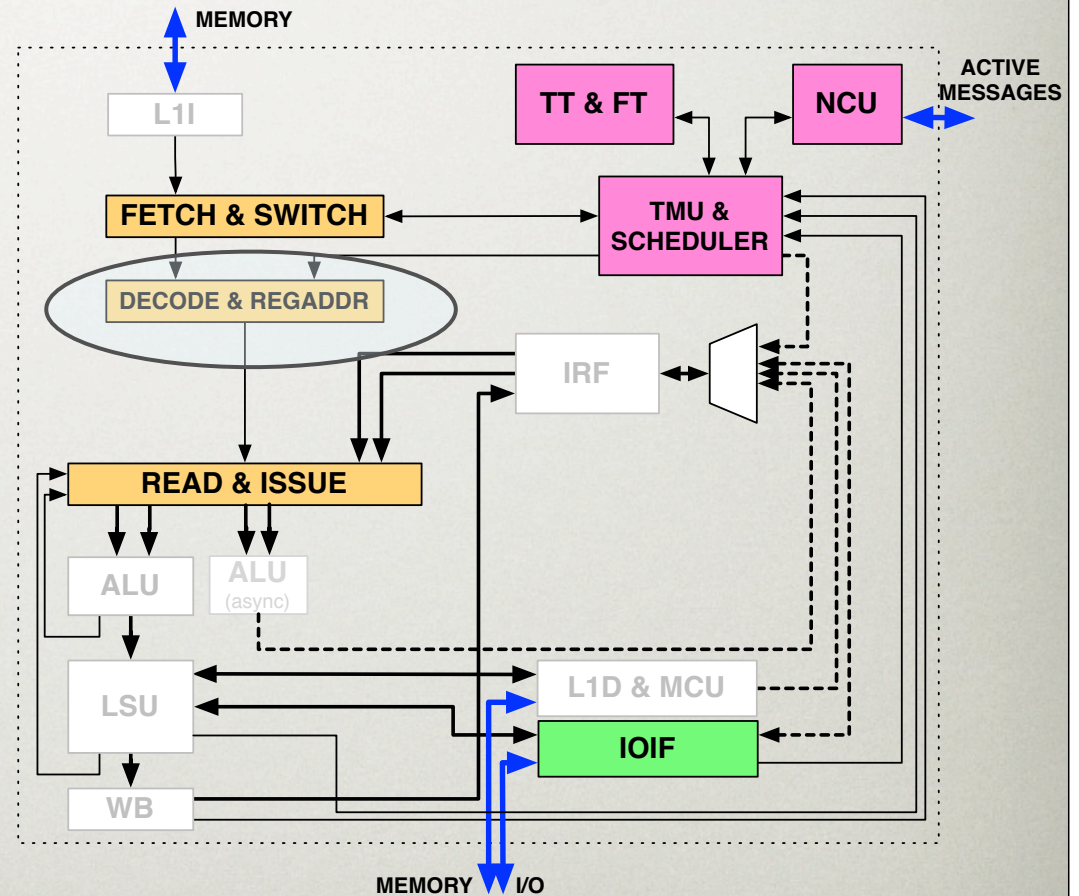


# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

Each thread has its own register window of configurable size

Base address in IRF computed during decode

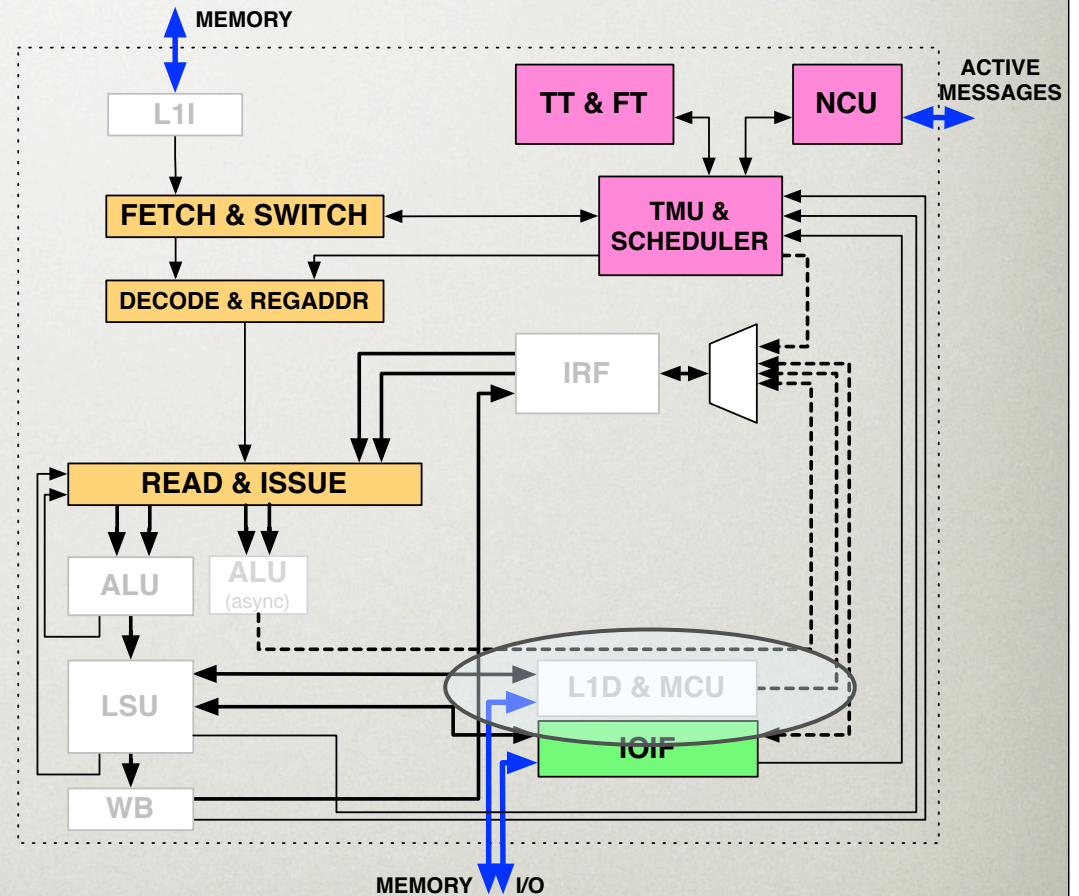




# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

Upon L1 miss the thread is allowed to go through a “waiting” state is written to the output register, only further instructions dependent on it will suspend





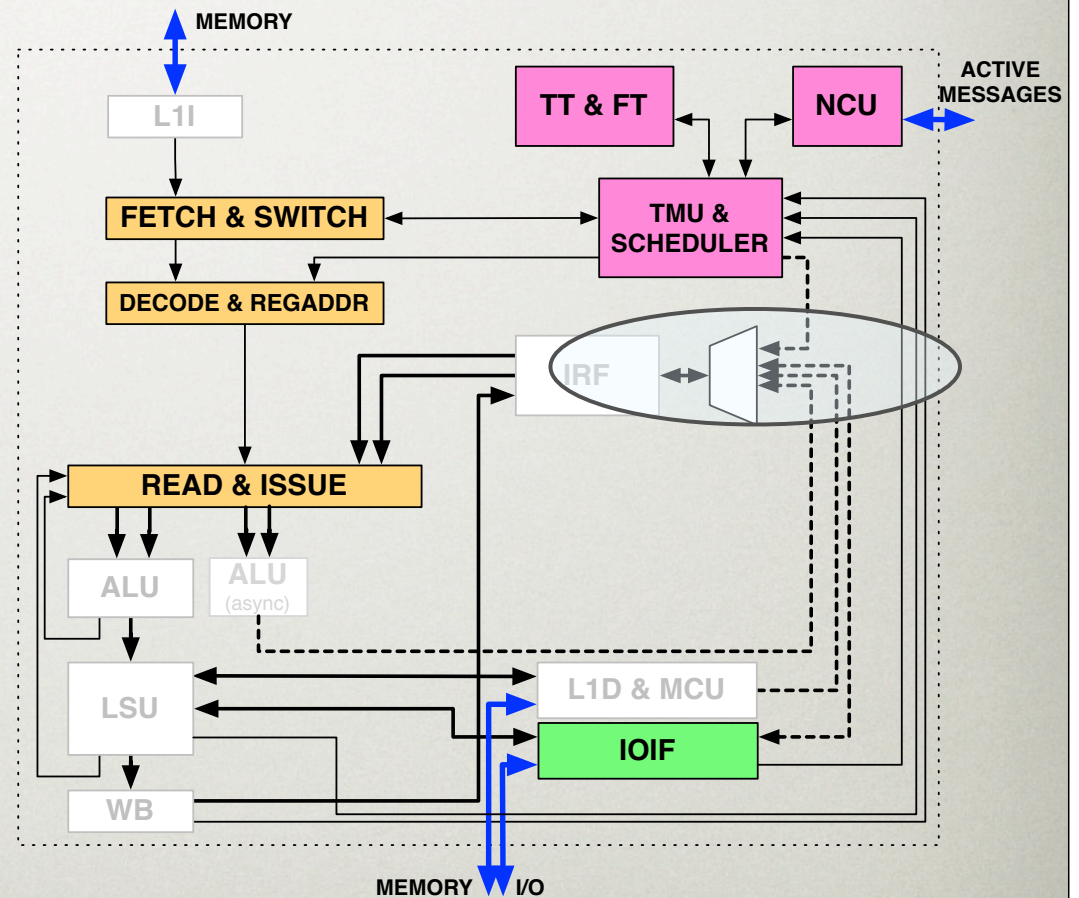
# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

Upon completion of a long-latency operation

(multi-cycle ALU, memory load)  
the value is written back  
asynchronously to the RF

this wakes up  
any waiting thread(s)  
by bringing them back  
on the active FIFO



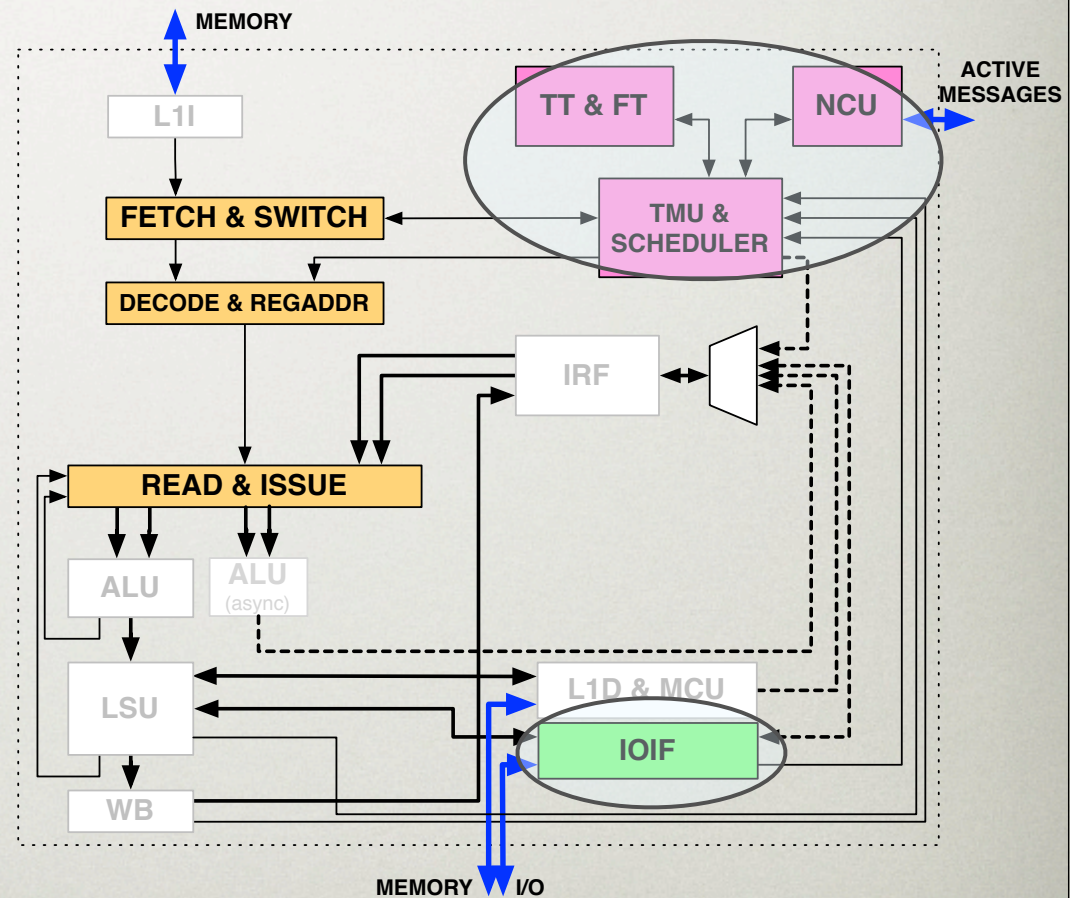


# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

The Thread Management Unit handles thread creation, synchronization, and provides thread metadata to the pipeline

Can be controlled from the pipeline (ISA extensions), control NoC messages or other I/O events





# APPLE-CORE

## A PERSPECTIVE SHIFT

<p>CORE I7 LINUX</p>	<p><b>Thread creation</b> (pre-allocated stack)</p> <p><b>&gt;10000 cycles in pipeline</b></p>	<p><b>Context switch</b></p> <p>syscalls, thread switch, trap, interrupt</p> <p><b>&gt;10000 cycles in pipeline</b></p>	<p><b>Thread cleanup</b></p> <p><b>&gt;10000 cycles in pipeline</b></p>
<p>D-RISC WITH TMU IN HARDWARE</p>	<p><b>Family creation</b> (metadata allocation for N threads)</p> <p><b>~15 cycles, async</b></p> <p><b>Thread creation</b> <b>1 cycle, async</b></p>	<p><b>Context switch</b></p> <p>at every waiting instruction, also I/O events</p> <p><b>&lt;1 cycles</b></p>	<p><b>Thread cleanup</b> <b>1 cycle, async</b></p> <p><b>Family cleanup</b> <b>2 cycles, async</b></p>



# APPLE-CORE

## A PERSPECTIVE SHIFT

---

<p>CORE I7</p>	<p><b>Function call</b></p> <p>with 4 registers spilled</p> <p><b>30-100 cycles</b></p>	<p><b>Predictable loop</b></p> <p>requires branch predictor + cache prefetching to maximize pipeline utilization</p> <p>2+ cycles / iteration overhead</p>
<p>D-RISC WITH TMU IN HARDWARE</p>	<p><b>Family creation</b></p> <p>of 1 thread, 31 "fresh" registers</p> <p><b>~15 cycles</b></p>	<p><b>Thread family</b></p> <p>1 thread / "iteration" reuses common TMU and pipeline no BP nor prefetch needed</p> <p>1 cycle / iteration overhead</p>



# APPLE-CORE

## A PERSPECTIVE SHIFT

---

	UNIT OF PROGRAMMING
IN A MOSTLY SEQUENTIAL CPU	<p>The sub-sequence, defined by structured <b>blocks</b></p> <p><b>Repetition</b> defined by loop control expression and <b>statement body</b></p>
IN A MULTITHREADED, MANY-CORE CHIP	<p>The <b>microthread</b> defined by <b>thread programs</b></p> <p><b>Repetition</b> defined by Family <b>logical index space</b> and <b>thread program</b></p>



# APPLE-CORE

## A PERSPECTIVE SHIFT

---

	SYNCHRONIZATION STYLE
IN A MOSTLY SEQUENTIAL CPU	<b>IMPLICIT</b>  One mechanism: <b>sequential ordering</b>
IN A MULTITHREADED, MANY-CORE CHIP	<b>SEMI-EXPLICIT</b>  <b>sequential</b> within threads <b>dataflow</b> within families <b>bulk synchronization</b> across families



# APPLE-CORE

## A PERSPECTIVE SHIFT

---

	BASIC MECHANISM FOR PROGRAM COMPOSITION
IN A MOSTLY SEQUENTIAL CPU	<p>Control flow transfer to/from <b>named entry point</b></p> <p>Defines <b>singular</b> functions and composition with <b>call/return</b> protocol</p> <p><b>Traps and interrupts</b> for async. events = implicit call/return</p>
IN A MULTITHREADED, MANY-CORE CHIP	<p><b>Family creation and family synchronization</b> of multiple threads running from the same <b>named entry point</b></p> <p>Defines <b>multi-way</b> functions and composition with <b>create/sync</b> protocol</p> <p>Optional: execution <b>placement</b> “at” a <b>named cluster of cores</b></p> <p>Asynchronous families for async. events = implicit create/sync</p>



# EVOLUTION OR REVOLUTION?

---

- Looks like a *revolution*:
  - **Can't reuse existing OS code as-is**
  - **Can't reuse existing low-level techniques**  
based on preemption and software schedulers  
eg. signals, interrupt handlers, "callbacks"
  - **Must use ISA concurrency in code generation**  
to exploit; requires language extensions  
and shakes **compiler assumptions**
- Can we really afford this?



# AN EVOLUTION, REALLY (1)

---

- Low-level machine code generation:
  - **Lift loop bodies as separate functions**
    - reuses techniques from GPU / accelerator world
  - **A thread is really a virtual processor**
    - threading is well-know in compilers already
  - **Higher-level compilers** can generate threaded low-level code from “productivity” languages
- Really a **convergence** of mature technology



# AN EVOLUTION, REALLY (2)

---

- **Managing asynchrony of “external” events**
  - I/O, traps, remote “syscalls”
- **An event handler is really a thread**
  - reuse as is, just entry / exit is different
- **Requires mutual exclusion of shared state**
  - already accepted in OS / library design
  
- **The benefit of extra bandwidth and lower latency will justify the req. adaptation, if any**

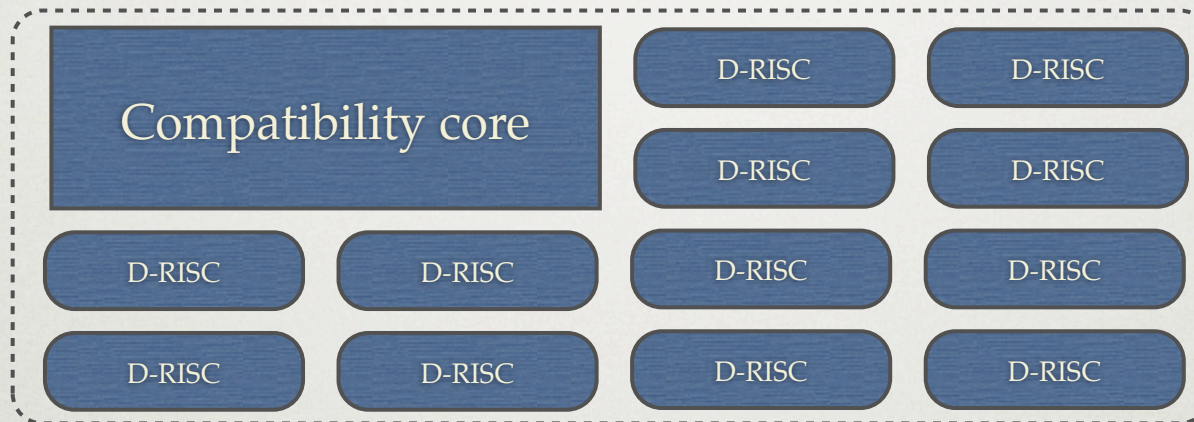


# AN EVOLUTION, REALLY

## (3)

---

- Sometimes **legacy OS and library code** cannot be adapted  
– typically device drivers, proprietary interfaces



- Solution: integrate a “compatibility” core on the same chip using **same NoC protocol for concurrency management**, then **delegate syscalls** behind library entry points
- With same ISA and shared memory **APIs can be kept as-is**
- The “**accelerator pattern**”, **transposed!** – similar to Cray XMT, on chip



# THE “REAL” ISSUES UNCOVERED IN APPLE-CORE

---

- **Validation:** how to detect detect errors, then compare with existing systems
  - need reference / base lines
- **Security:** isolation and observation
  - and how to reduce interference with performance
- **Resource management:**  
cores, but also **memory** and **NoC channels**
  - how to reduce management overheads
- NB: these issues are general to many-core, but they are exarcerbated in Apple-CORE



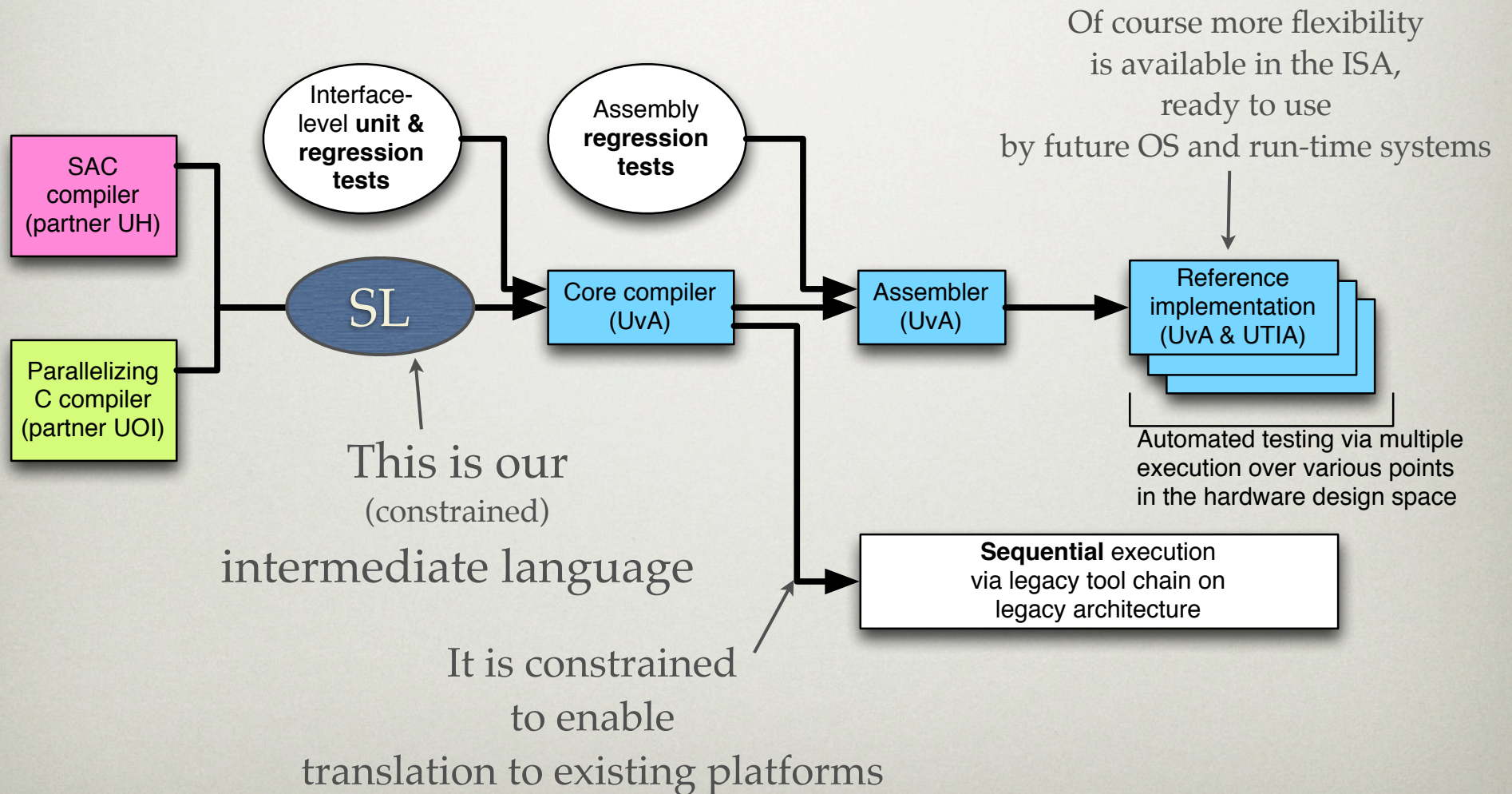
# VALIDATION

---

- Solution:
  1. Choose a **subset of the ISA** that can be emulated in legacy platforms
  2. Design the intermediate language SL to use only this subset to **constrain programs**
  3. Implement **compilation to both** the new platform and legacy systems and perform **comparative testing**
- This subset resembles **fork/join with families and forward-only dataflow synchronization**
- It is **deadlock-free**, mostly **deterministic** and **can be serialized** (cf Cilk, Chapel)



# VALIDATION





# SECURITY

---

- Isolation via virtual address spaces is expensive
- Instead Apple-CORE has fine-grain **capabilities** on memory and concurrency management  
Concept imported from microkernels: Opal, Mungi, etc.  
VT still available at the level of core cluster or entire chip
- **Observation** via **out-of-band hw monitors**  
Can observe performance **remotely**  
**including TMU** – active/suspended threads, etc.  
“Out of band” means in parallel to pipeline execution, also potentially **separate from memory** on different NoC



# RESOURCE MANAGEMENT

---

- At the finest grain:  
provide TLS to threads created by TMU  
Solution: **pre-allocate** and **partition statically**
- **Algorithms:** distributed memory allocator, garbage collection using reference counting
- **Application components:**  
OS allocates and deallocates cores, memory and network links for top-level family entry points  
– this is called **SEP** and is **distributed**
- Either **explicit allocation** in programs (Apple-CORE)  
Or annotated static requirements, aggregated at run-time by RTS/OS (upcoming project ADVANCE)



# THE SYSTEM

## VIRTUALIZATION PLATFORM

---

- Apple-CORE: was bottom-up from DRISC
- Main outcome: common **concurrency management protocol** on chip
- Can be generalized to other many-core chips: define **placed family creation** as basic primitive instead of function calls, message passing and IPIs
- This defines a **new form of processor virtualization** where **concurrency is the norm**
- Apple-CORE makes this most efficient (in hw) but the protocol **can be used with other platforms**



# FOREGROUND PRODUCED

---

- Technology:
  - various **simulators** for a many-core chip with hardware concurrency management (Microgrid)
  - **MGOS: OS and library components** to drive the hardware architecture, including **resource allocators** and **API bridges**
  - **compilation tools** to/from the SVP intermediate language
  - **software run-time systems** for commodity multi-cores using SVP semantics
  - **Tests and benchmarks** to validate and evaluate fine-grained concurrency management
- + Accompanying documentation & know-how



# FOREGROUND PRODUCED

---

Common C language primitives

MGOS

Hydra

ptl

LPEL

HLSim

MGSim

UTLEON3

Distributed memory  
hw multithreaded  
multi-cores

Shared memory  
sw multithreaded  
multi-cores

Microgrid hardware model



# SUMMARY

---

- A true **perspective shift** for the basic OS/ compiler abstractions:
  - from *sequence* to *concurrency*
  - from *loops* to *microthreads*
  - from *function calls* to *family creation*
  - new focus on **placement** and locality
- Revolution in hardware, yet only an evolution in software
- Middle ground: a common set of primitives  
= a **concurrency management protocol** on chip
- This is generalized from D-RISC towards portable SVP



# THANK YOU!

---

- More information:
  - <http://www.apple-core.info/>
  - <http://www.svp-home.org/>



# CONCURRENCY MANAGEMENT PROTOCOL

---

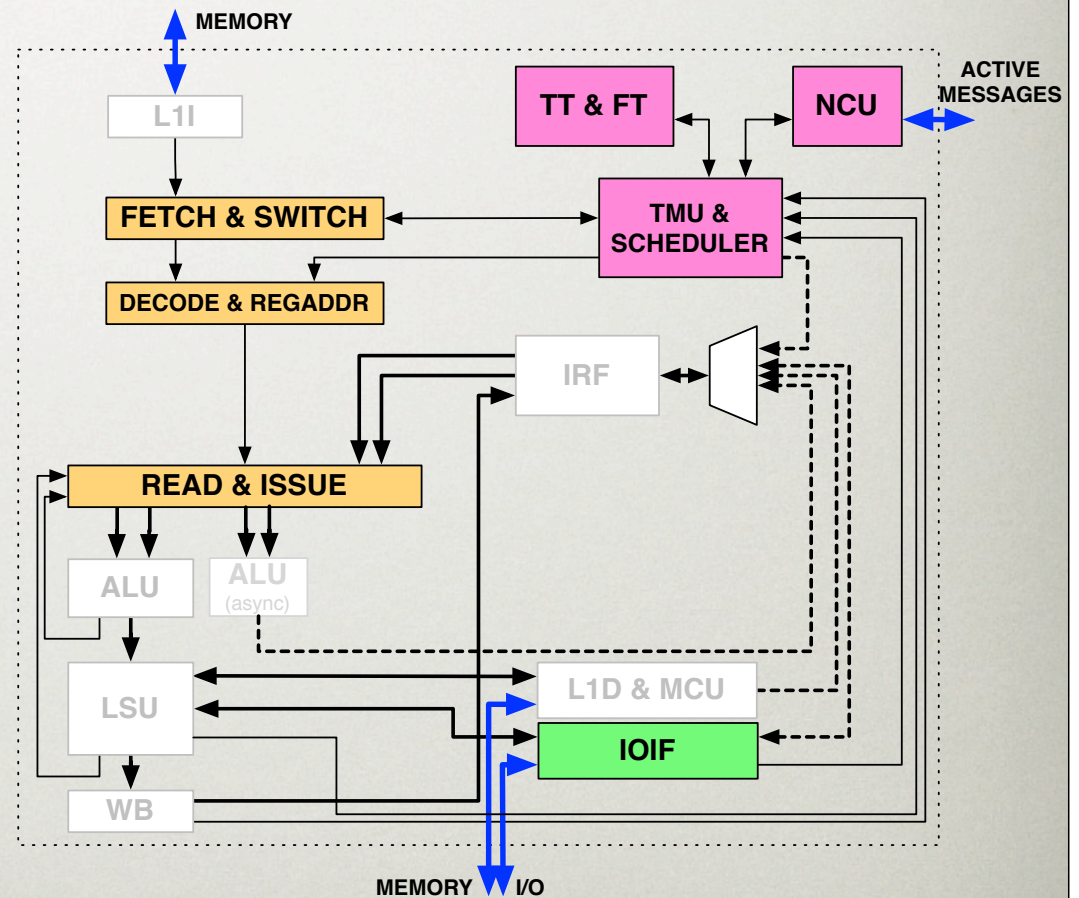
allocate \$Place $\rightarrow$ \$F	Allocate a family context
setstart / setlimit / setstep / setblock \$F, \$V $\rightarrow$ $\emptyset$	Prepare family creation
create \$F, \$PC $\rightarrow$ \$ack	Start bulk creation of threads
rput \$F, R, \$V $\rightarrow$ $\emptyset$ rget \$F, R $\rightarrow$ \$V	Read / write dataflow channels remotely
sync \$F $\rightarrow$ \$ack	Bulk synchronize on termination
release \$F $\rightarrow$ $\emptyset$	De-allocate a family context



# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

- Reference parameters:
  - 1GHz clock
  - 1Kregs IRF  
+ 1Kregs FRF
  - 256 thread contexts  
32 family contexts
  - 64-bit  
Alpha-like ISA
  - Asynchronous FPU,  
2 cores / FPU





# THE D-RISC CORE

## OPERATING SYSTEM IN HARDWARE

- On FPGA (UTLEON3)
  - 20MHz clock
  - 512-1024 registers
  - 64-128 thread contexts
  - 16-32 family contexts
- 32-bit SPARC-like ISA

