

# **OPERATING SYSTEMS IN HARDWARE**

## **SCALING FROM 10 TO 1000 CORES**

**RAPHAEL 'KENA' POSS**  
**UNIVERSITY OF AMSTERDAM, THE NETHERLANDS**

**IN COOPERATION WITH**  
**CHRIS JESSHOPE, MIKE LANKAMP, MICHEL VAN TOL, AND MANY OTHERS**  
**AT THE CSA GROUP**



# CURRENT GENERAL-PURPOSE MULTI-CORES ARE BASED ON LEGACY

---

- Historical focus on single-thread performance  
(developments in general-purpose processors: registers, branch prediction, prefetching, out-of-order execution, superscalar issue, trace caches, etc.)
- Legacy heavily **biased towards single threads:**
  - Symptom: **interrupts** are the **only way** to signal asynchronous external events
  - Retro-fitting **hardware multithreading** is **difficult** because of the sequential core's complexity
- **What if...**  
**we redesigned general-purpose processors,**  
**assuming concurrency is the norm in software?**

# THE APPLE-CORE APPROACH

---

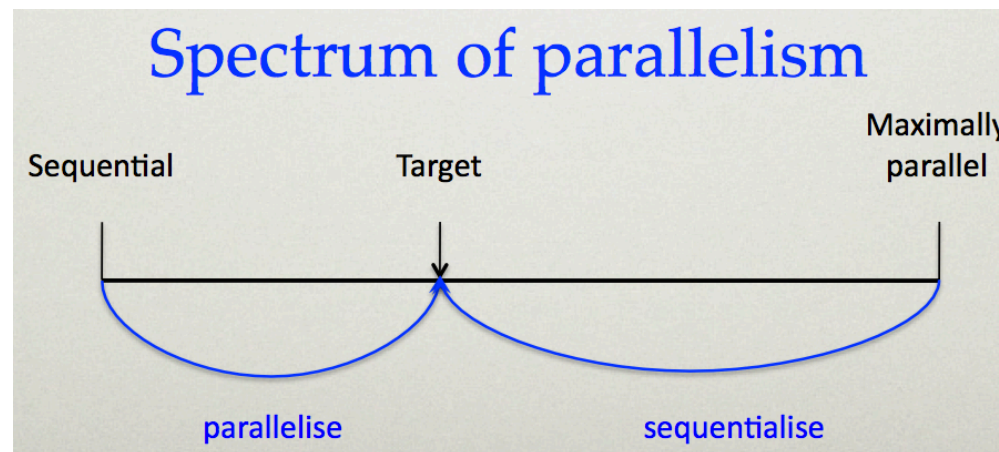
- Hardware threading is not new... Denelcor HEP introduced this 30 years ago
- Apple-CORE's contribution is in the *system interface with a consistent view on concurrency management* whatever resources are available
  - from a single thread to many cores
- The goal is to *make concurrency the norm* and the processing resource fungible
  - like money it must be divisible and interchangeable

# THE APPLE-CORE APPROACH

---

- With massive concurrency predicted, the Apple-CORE approach is to:
  - capture maximal concurrency in the binary interface
  - automatically sequentialise it based on resource availability at run time

Past  
approach



Apple-CORE  
approach

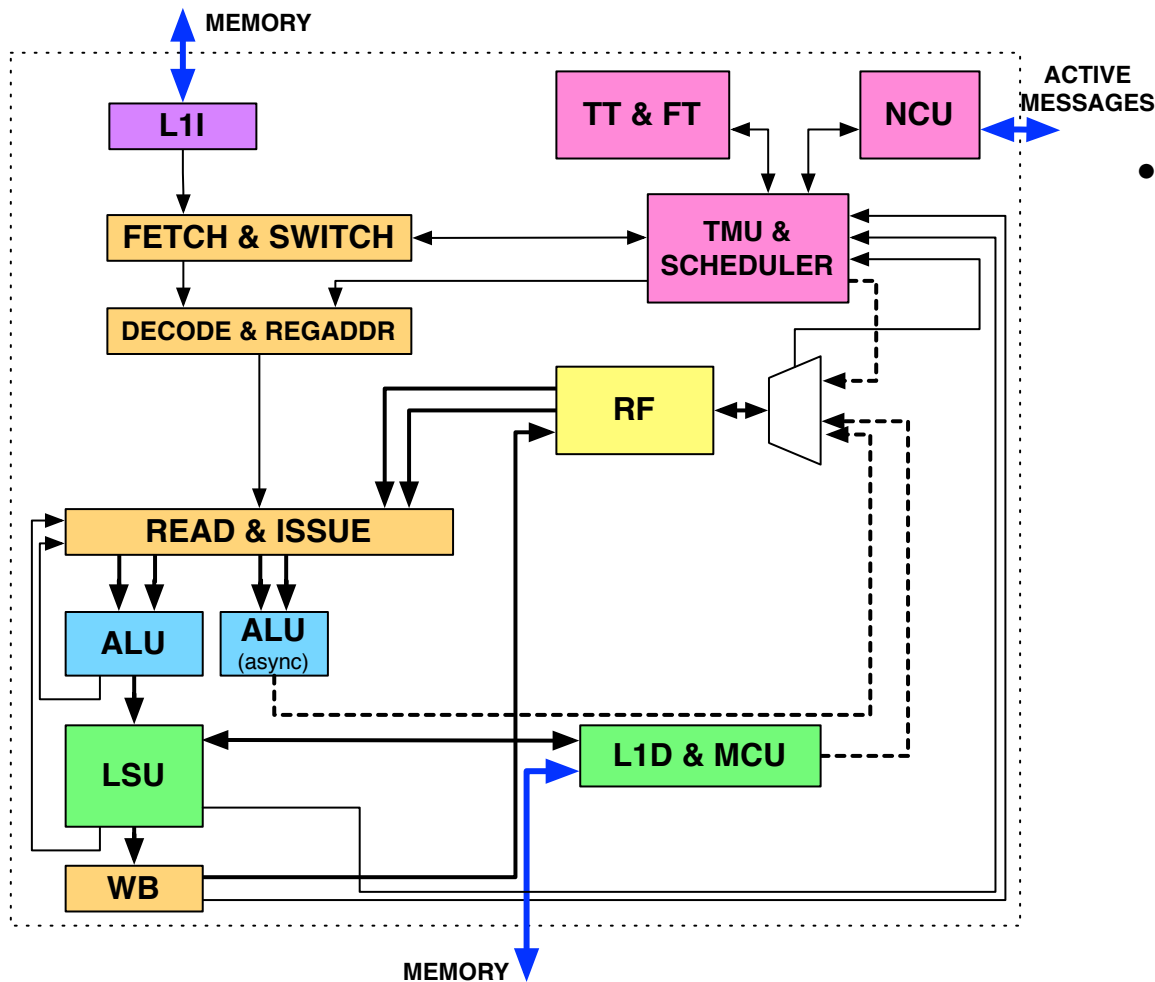
# KEY ABSTRACTIONS

---

- Dataflow synchronisation on instruction execution - *latency tolerance*
- Deterministic programming model based on multi-way fork / sync - *create families of threads*
- Memory consistent at fork and sync events
- Handles on resources - *places*
- Asynchronously execution of *families* at *places*
- Capture sequence where necessary by thread- to- thread dependencies - *shared variables*

# THE D-RISC CORE

# D-RISC CORES

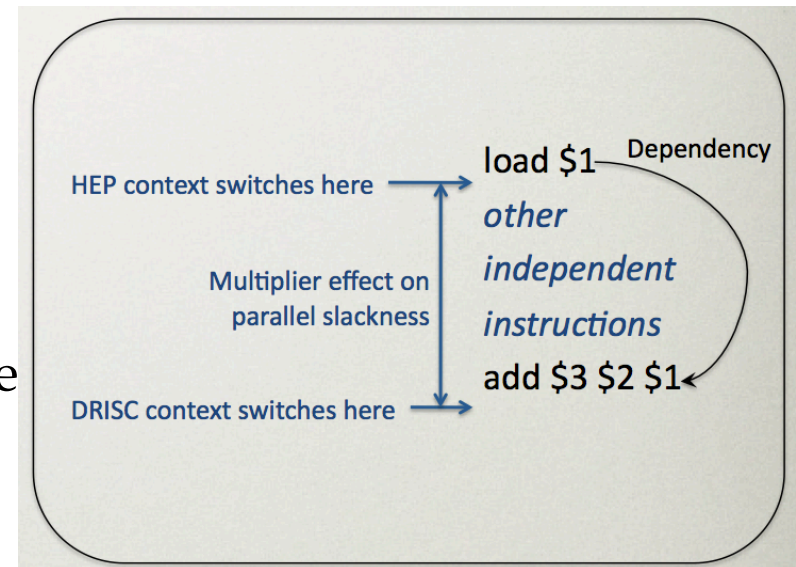


- D-RISC cores:  
hardware multithreading +  
dynamic dataflow scheduling
- fine-grained threads: 0-cycle  
thread switching, <2 cycles  
creation overhead
- ISA instructions and NoC  
protocol for thread management
- dedicated hardware processes  
for bulk creation and  
synchronization
- No preemption/interrupts;  
events “create” new threads

In-order, single-issue RISC: small, cheaper, faster/watt

# DATAFLOW SCHEDULING

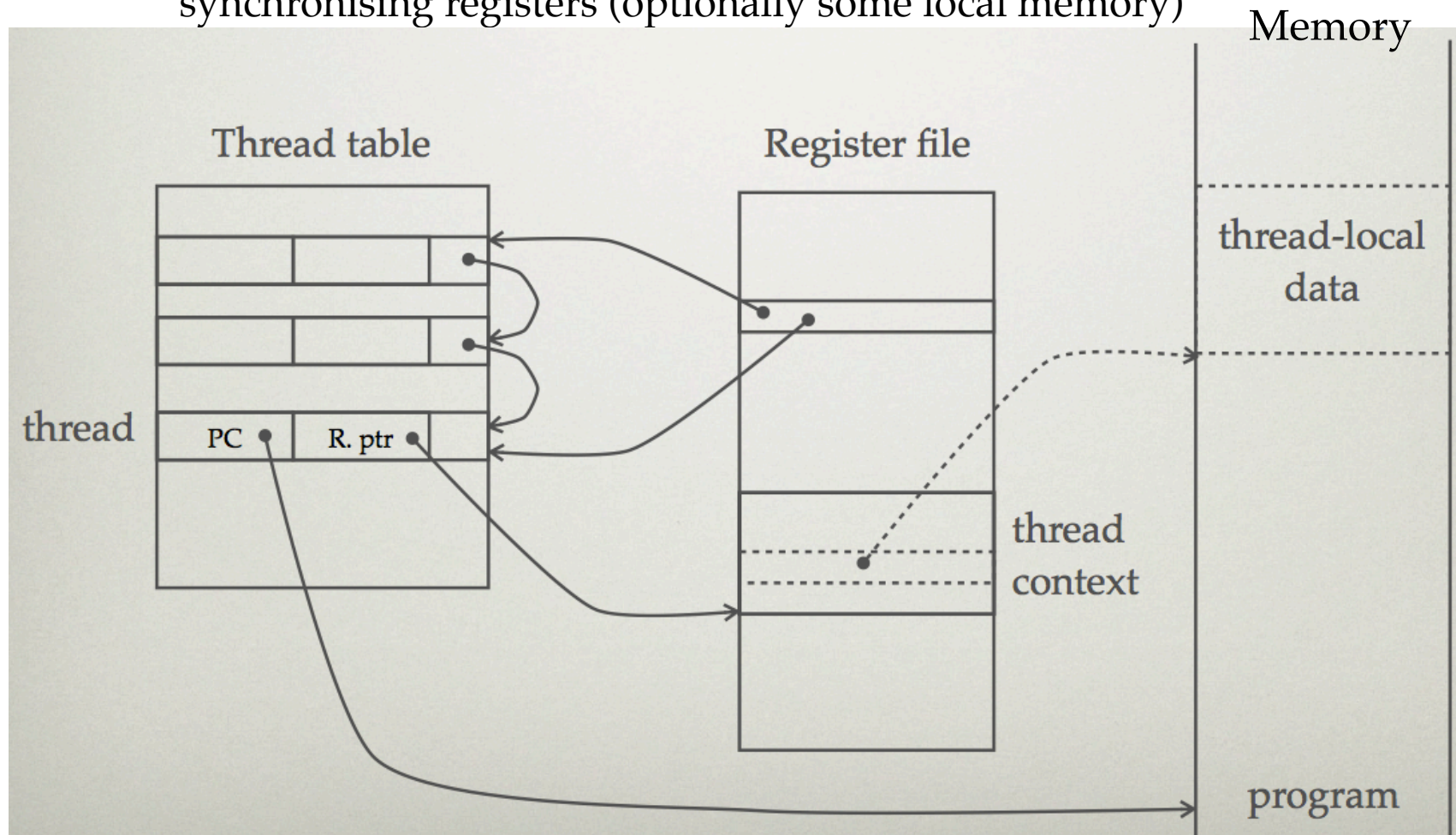
- Dataflow is a more effective way to tolerate latency in instruction execution - e.g. memory accesses take 100s or even 1000s of cycles
- Prior approaches e.g. HEP/ Niagara context switch on the hazard-inducing instruction
- DRISC synchronises on the result of the hazard, i.e. captures the dependency - *dataflow scheduling*
- This requires a synchronising register file (i-structures) and an efficient mechanism to store and manage thread continuations





# D-RISC CONTINUATIONS

- Continuation comprise a PC and a context comprising from 2 to 32 synchronising registers (optionally some local memory)



# ASYNCHRONY IN D-RISC CORES

---

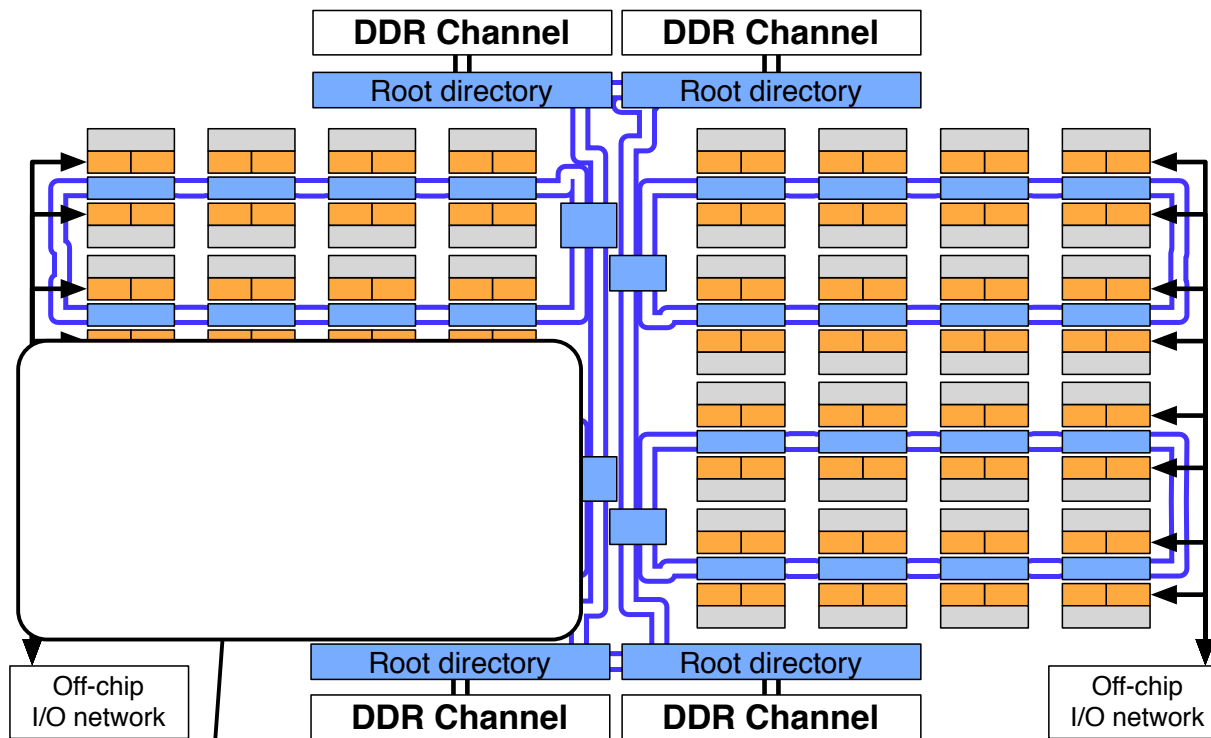
- *Intra-thread* - i.e. an instruction that is issued and synchronised within the same thread:
  - Memory references
  - FPU operations (NB: shared FPU)
  - Functions implemented as logic
  - Create a family and Sync on termination of a family
- *Inter-thread* - read / write to a remote register file
  - This asynchronously activates a waiting thread and is one of two mechanism used for system signalling - the other is asynchronously creating a family of thread(s)

# ACHIEVEMENTS - D-RISC CORES

---

- *Single DRISC core - FPGA*  
implementation of single core  
(SPARC ISA)
- *Many-core Microgrid*  
(Alpha and SPARC ISAs)
- *Low-level C-based compiler, operating  
system and run-time environment for both  
platforms*

# EXAMPLE 128-CORE MICROGRID



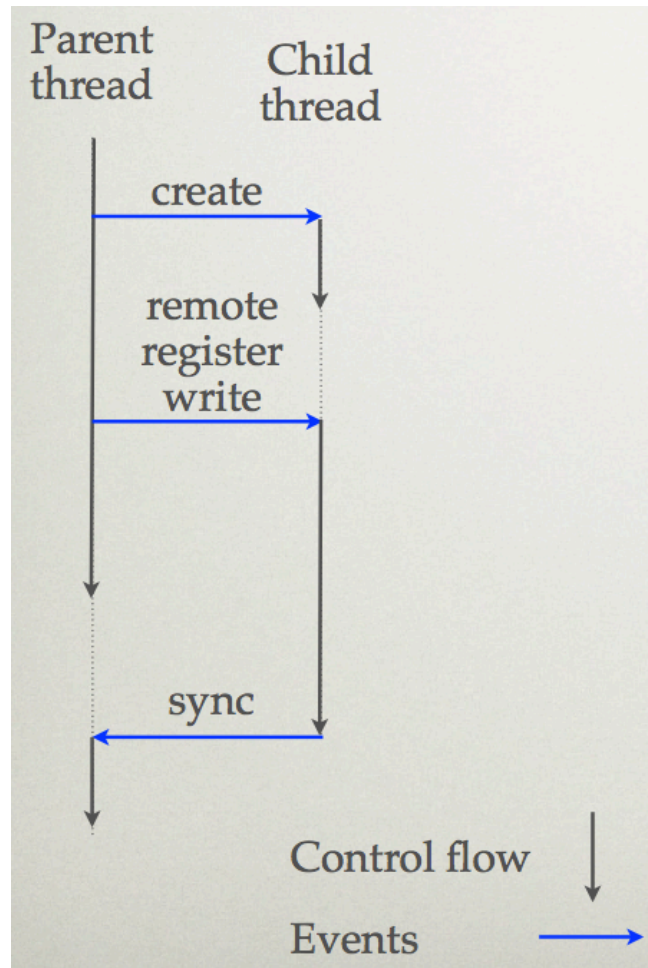
Approximate size of one Nehalem (i7) core  
for comparison

- 32000+ hw threads
- 5MB distributed cache
- shared MMU  
= single virtual address space, protection using capabilities
- Weak cache coherency
- no support for global memory atomics – instead synchronization using point-to-point messaging

Area estimates with CACTI: 100mm<sup>2</sup> @ 35nm

# **MICROTHREADED PROGRAMMING**

# FUNCTIONS AS THREADS

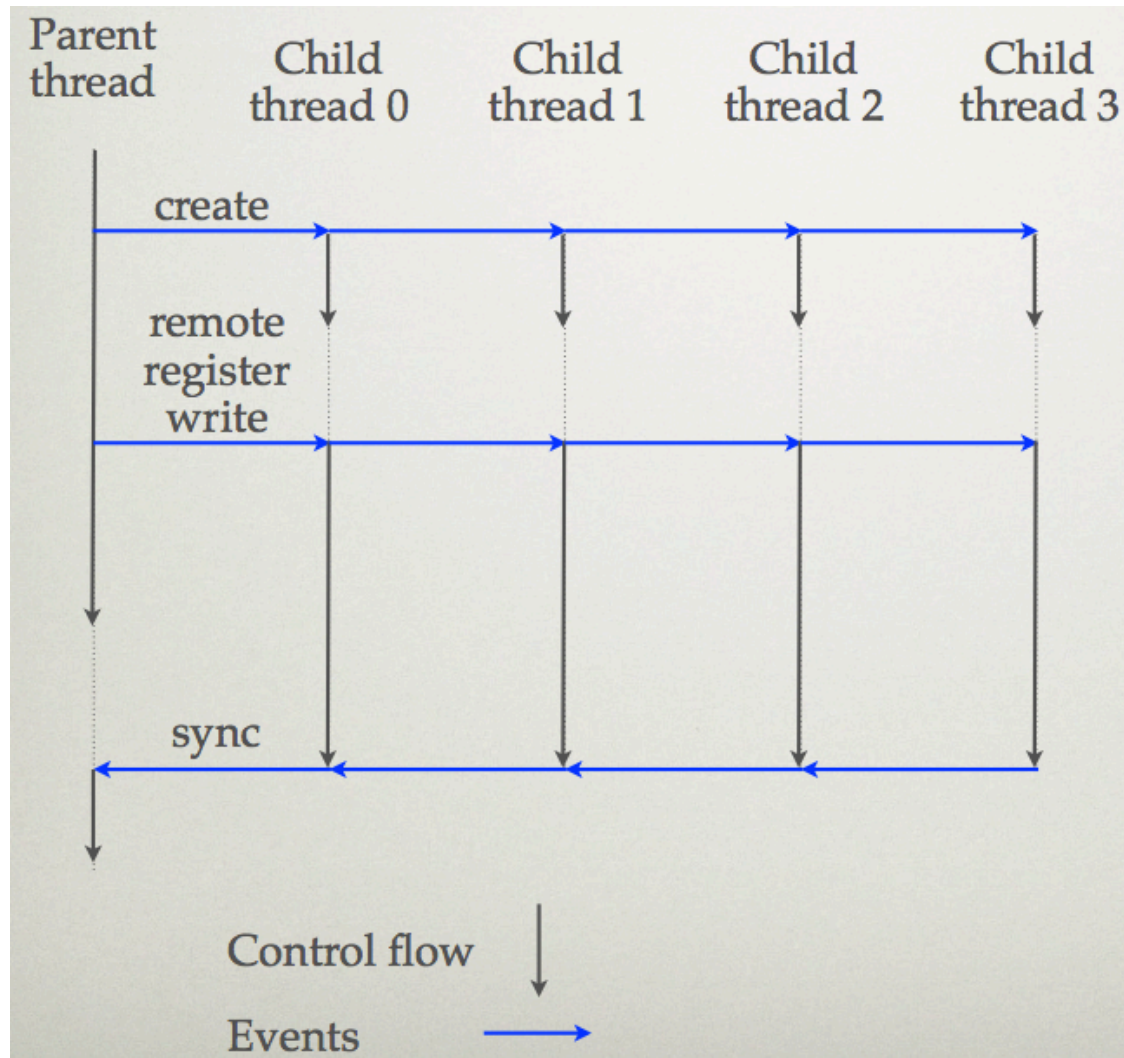


- Events are **create/sync** and **remote register read/writes**
- register writes are synchronising
- asynchronous parameter passing
- create executed early and parent and child execute concurrently

```
foo(A)
{...
...A...
}...
```

```
...
create() foo(P)
... P:=... ...
sync
```

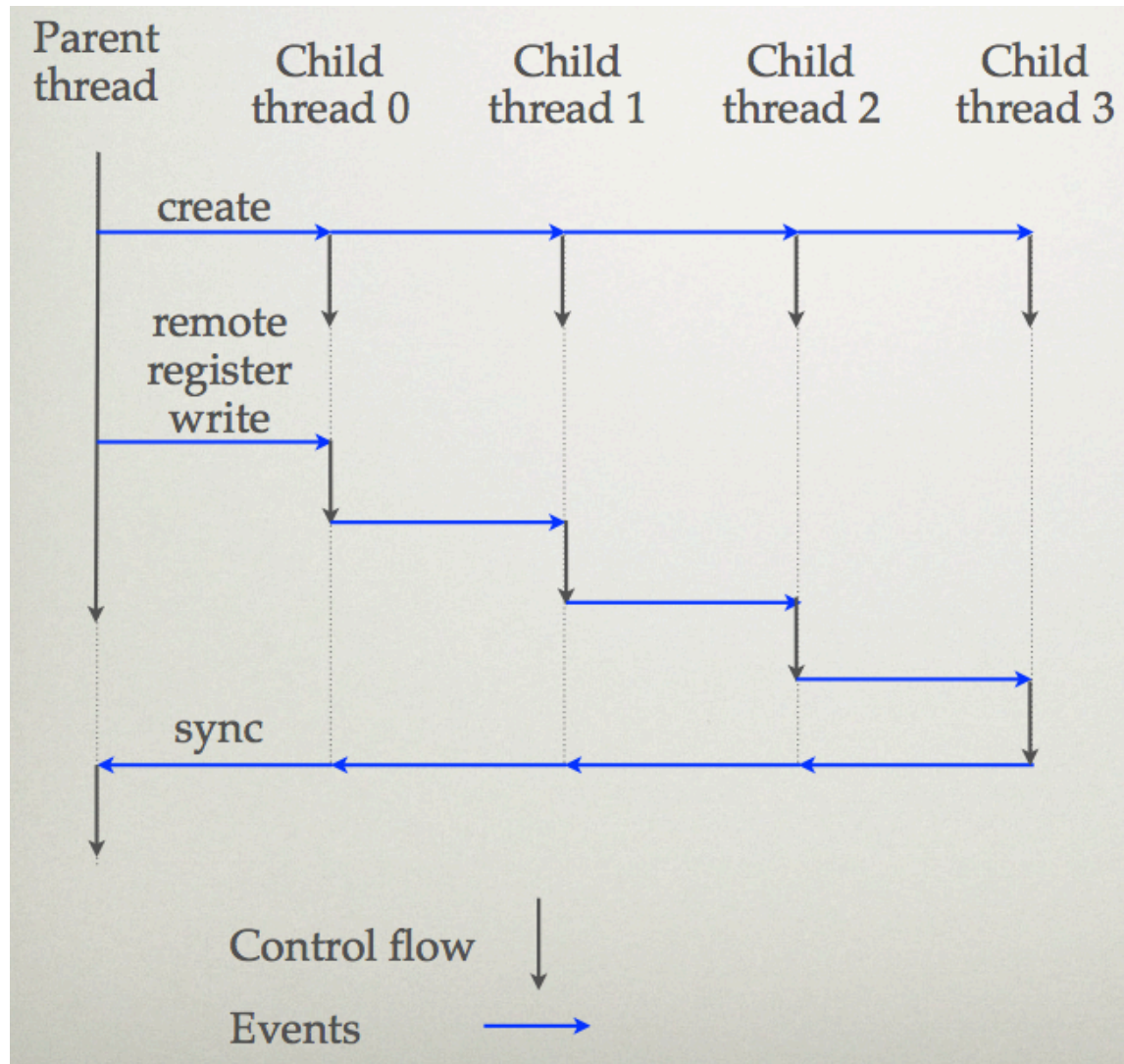
# LOOPS AS THREADS



```
foo(A)
{...
...A...
}...
```

```
...
create(i=0..3)
  foo(P)
... P:=... ...
sync
```

# SEQUENCE AS THREADS



```
foo(A)
```

```
{...
```

```
  A := ...A...
```

```
}...
```

```
...
```

```
create(i=0..3)
```

```
  foo(P)
```

```
... P:=... ..
```

```
sync
```



# WHY CAPTURE SEQUENCE AS THREADS?

---

- Still need local concurrency for latency tolerance.  
*e.g. naive inner product*
  - $\text{sum} := \text{sum} + A[i] * B[i]$
  - compiles to a six instruction thread in Alpha
  - $A+i$ ,  $B+i$ , `load`, `load`, `fmul` are all independent
  - only one instruction: `fadd $S1 $D1 $L1` is constrained to be executed sequentially
- No speedup, but *independent instructions can be scheduled concurrently* on a single core
  - This allows multiple concurrent memory loads and floating point multiplies to be executed concurrently on a single core

# PROGRAMMING ENVIRONMENT

---

- C with extensions -  $\mu$ TC for documentation, SL in implementations
- New *language primitives* for concurrency creation, synchronization
- A *thread program* is like a C function with special constructs for dataflow channels
- *No model asymmetry* (like with CUDA / OpenCL): any thread program can use any library service
- To be *targeted by higher-level compilers*

# 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 utilization</p> <p>1+ cycles / iteration overhead</p>
<p>D-RISC WITH TMU IN HARDWARE</p>	<p><b>Bulk thread creation</b></p> <p>of 1 thread, 31 "fresh" registers</p> <p><b>~15 cycles</b> (7c sync, ~8c async)</p>	<p><b>Thread family</b></p> <p>1 thread / "iteration" reuses common TMU and pipeline no BP nor prefetch needed</p> <p>0+ iteration overhead</p>

# **OPERATING SYSTEMS FOR MICROGRIDS**

# EVOLUTION OR REVOLUTION?

---

- Looks like a *revolution*:
  - **Can't (shouldn'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)

---

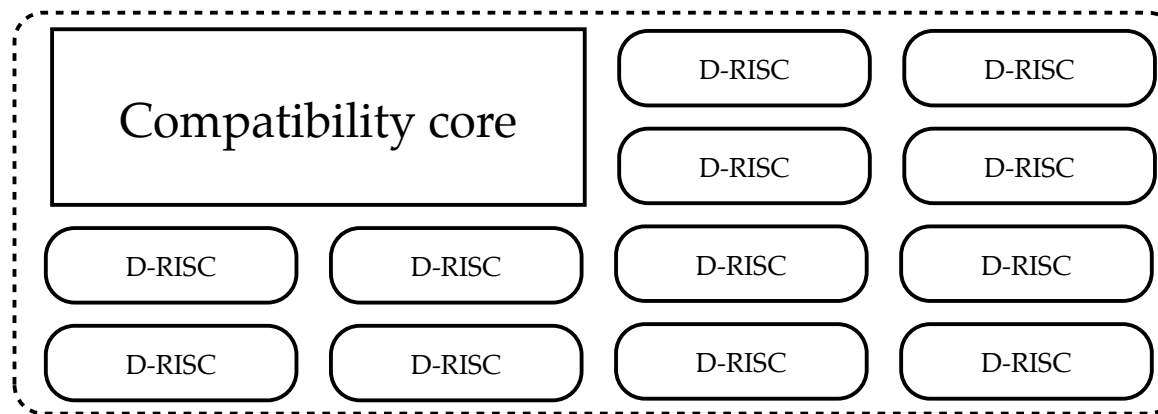
- Managing asynchrony of “external” events
  - I/O, traps, remote “syscalls”
- *An event handler is really a thread*
  - reuse as is, just entry / exit is different
- Requires *mutual exclusion of shared state*
  - already accepted in OS / library design
- The benefit of extra bandwidth and lower latency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



# FOREGROUND PRODUCED

---

- Technology:
  - various *emulators/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

# OS 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 in ISA  
= *a concurrency management protocol* on chip
- This is generalized from D-RISC towards portable SVP

# THE “MAIN” ISSUES UNCOVERED IN APPLE-CORE

---

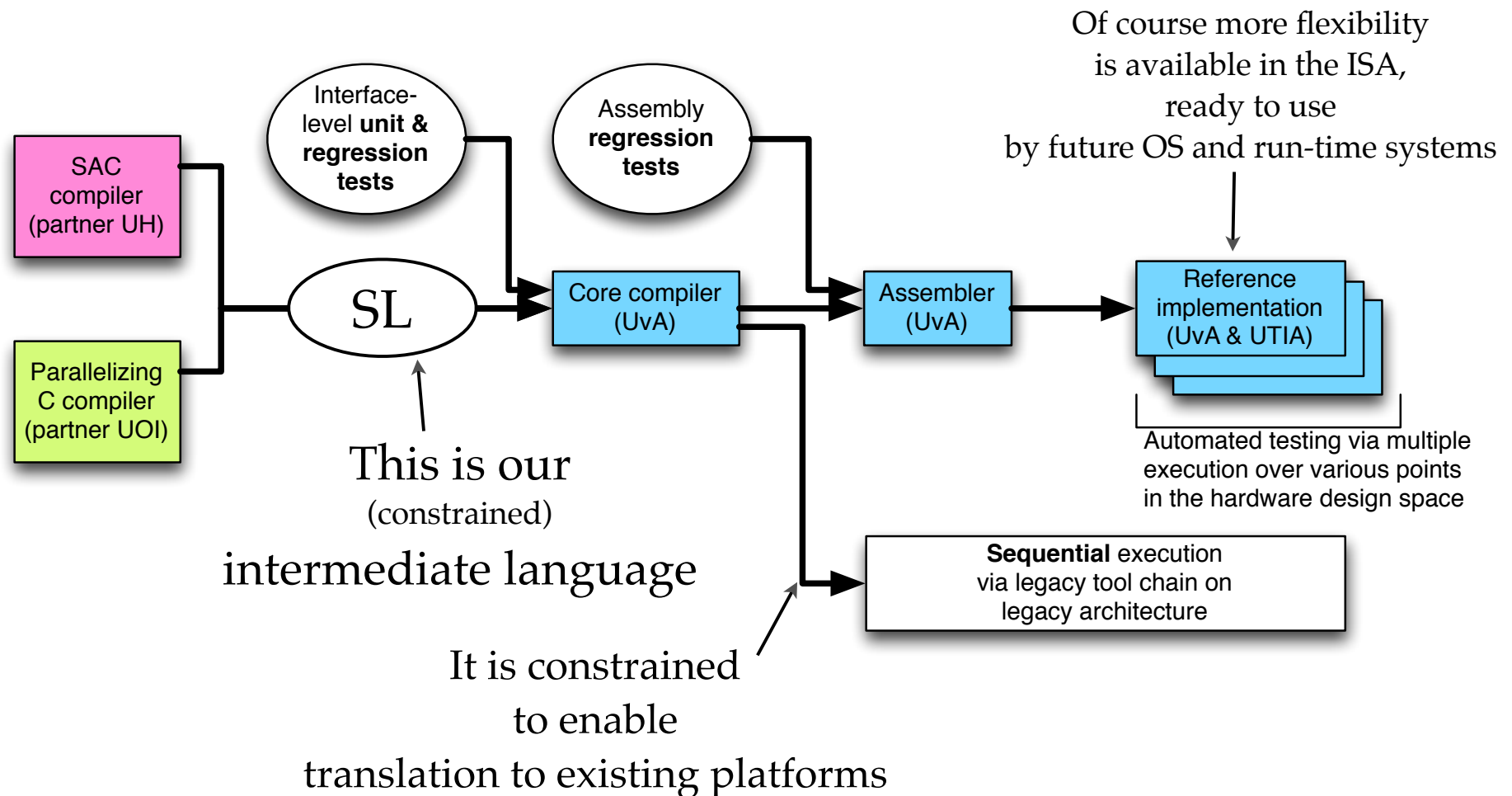
- **Validation:** how to detect detect errors, then compare with existing systems
  - need reference / base lines
- **Resource management:**  
**cores, but also memory and NoC channels**
  - how to reduce management overheads
- NB: these issues are general to all many-core processors, but exarcerbated 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, deterministic* if race-free and *can be serialized* (cf Cilk, Chapel)

# VALIDATION



# RESOURCE MANAGEMENT

---

- At the finest grain:  
*provide TLS to threads* created by TMU  
Solution: *pre-allocate* and *partition statically*
- *Concurrency resources*: let programs define more concurrency than available, serialize on demand
- *Algorithms*: distributed memory allocator, garbage collection using reference counting

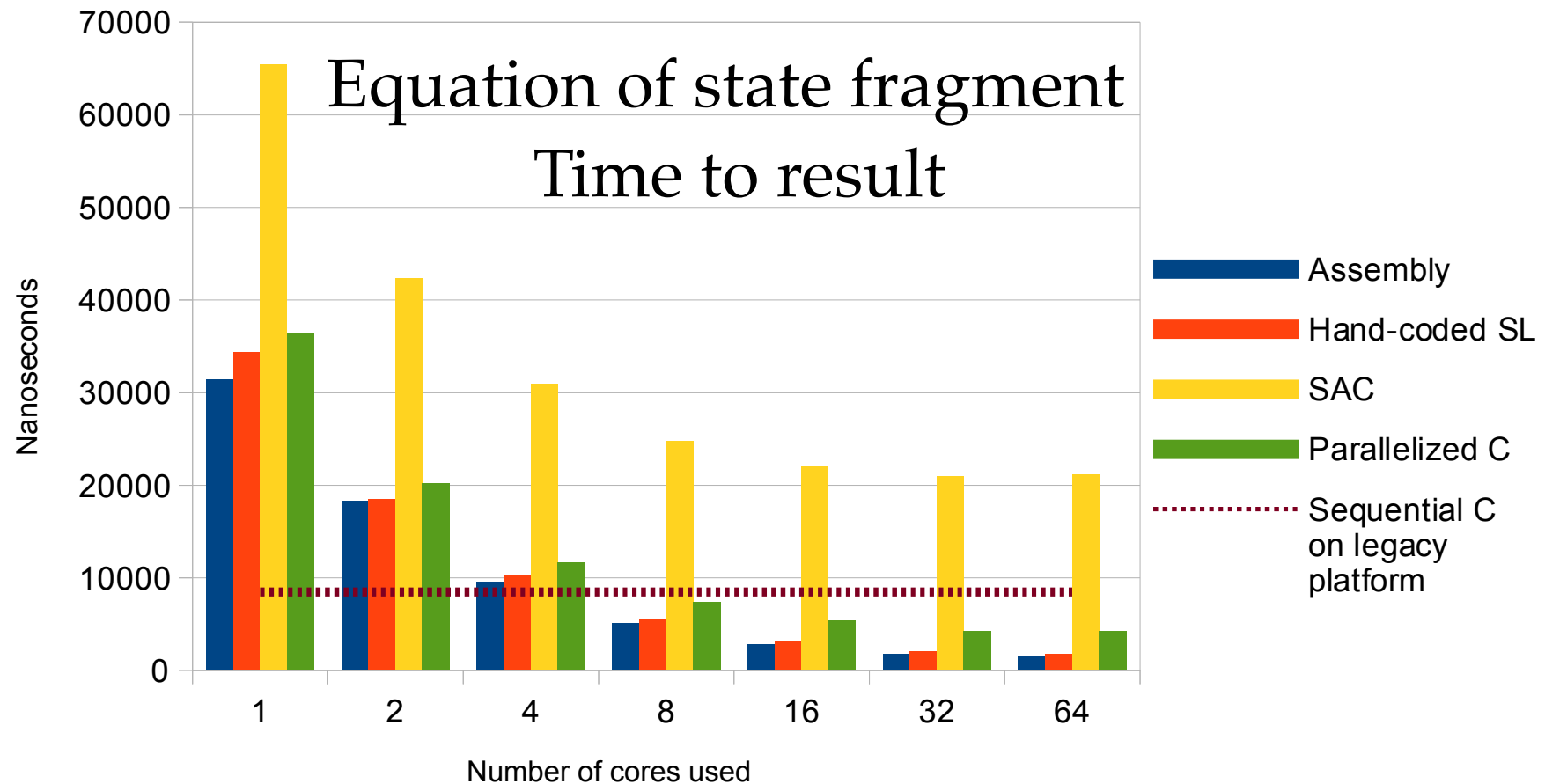
# RESOURCE MANAGEMENT

---

- *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  
  
Or *annotated static requirements*, aggregated  
at run-time by RTS / OS

# RESULTS:

## MEMORY-BOUND KERNELS

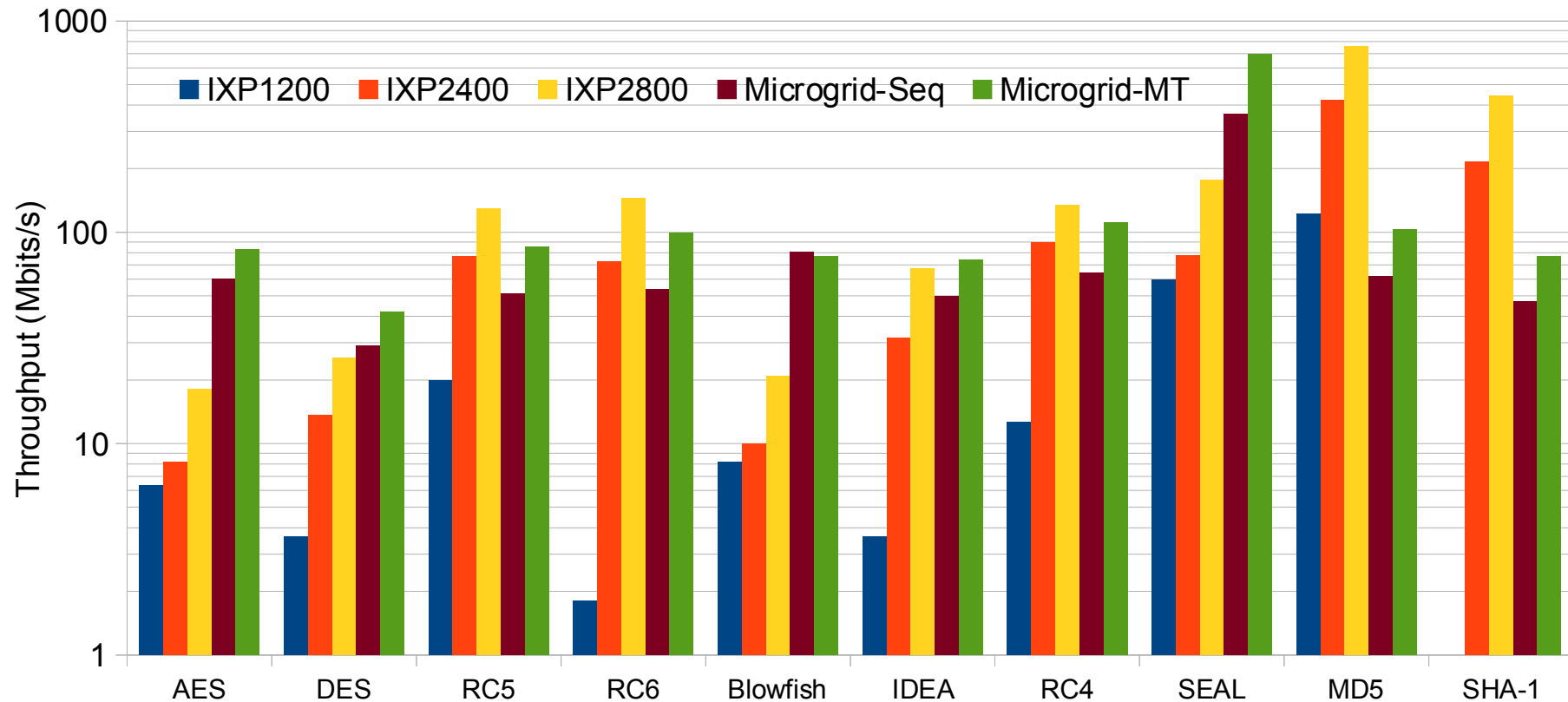


Legacy platform = MacBook Pro, Core 2 Duo @ 2.4GHz

area(1 Core 2 Duo core) ~ area(32 Microgrid cores)



# RESULTS: THROUGHPUT WORKLOADS



Intel IXP = embedded processor specialized for cryptographic workloads

Main results: **Microgrids are general-purpose**, ie not specialized  
**yet compete** on throughput with state-of-the art specialized hardware

# RESULTS, WHAT'S NEXT?

---

- *Internal* issues: memory consistency, scalable cache protocols, ISA semantics, etc.
- *External* issues from outside architecture: platform virtualization, space scheduling, I/O device drivers
- *Fundamental* issues: concurrent complexity theory

# THANK YOU!

---

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



# SVP CONCURRENCY MANAGEMENT PROTOCOL

---

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

# EXTRA - A PERSPECTIVE SHIFT

---

<p>CORE I7 LINUX</p>	<p><b>Thread creation</b>  (pre-allocated stack)  <b>&gt;10000 cycles in pipeline</b></p>	<p><b>Context switch</b>  syscalls, thread switch, trap, interrupt  <b>&gt;10000 cycles in pipeline</b></p>	<p><b>Thread cleanup</b>  <b>&gt;10000 cycles in pipeline</b></p>
<p>D-RISC WITH TMU IN HARDWARE</p>	<p><b>Bulk creation</b> (metadata allocation for N threads) <b>~15 cycles</b> (7c sync, ~8c async)  <b>Thread creation</b> <b>1 cycle, async</b></p>	<p><b>Context switch</b>  at every waiting instruction, also I/O events  <b>&lt;1 cycles</b></p>	<p><b>Thread cleanup</b> <b>1 cycle, async</b>  <b>Bulk synchronizer cleanup</b> <b>2 cycles, async</b></p>

# FOREGROUND PRODUCED

---

Common C language primitives

MGOS

Hydra

ptl

HLSim

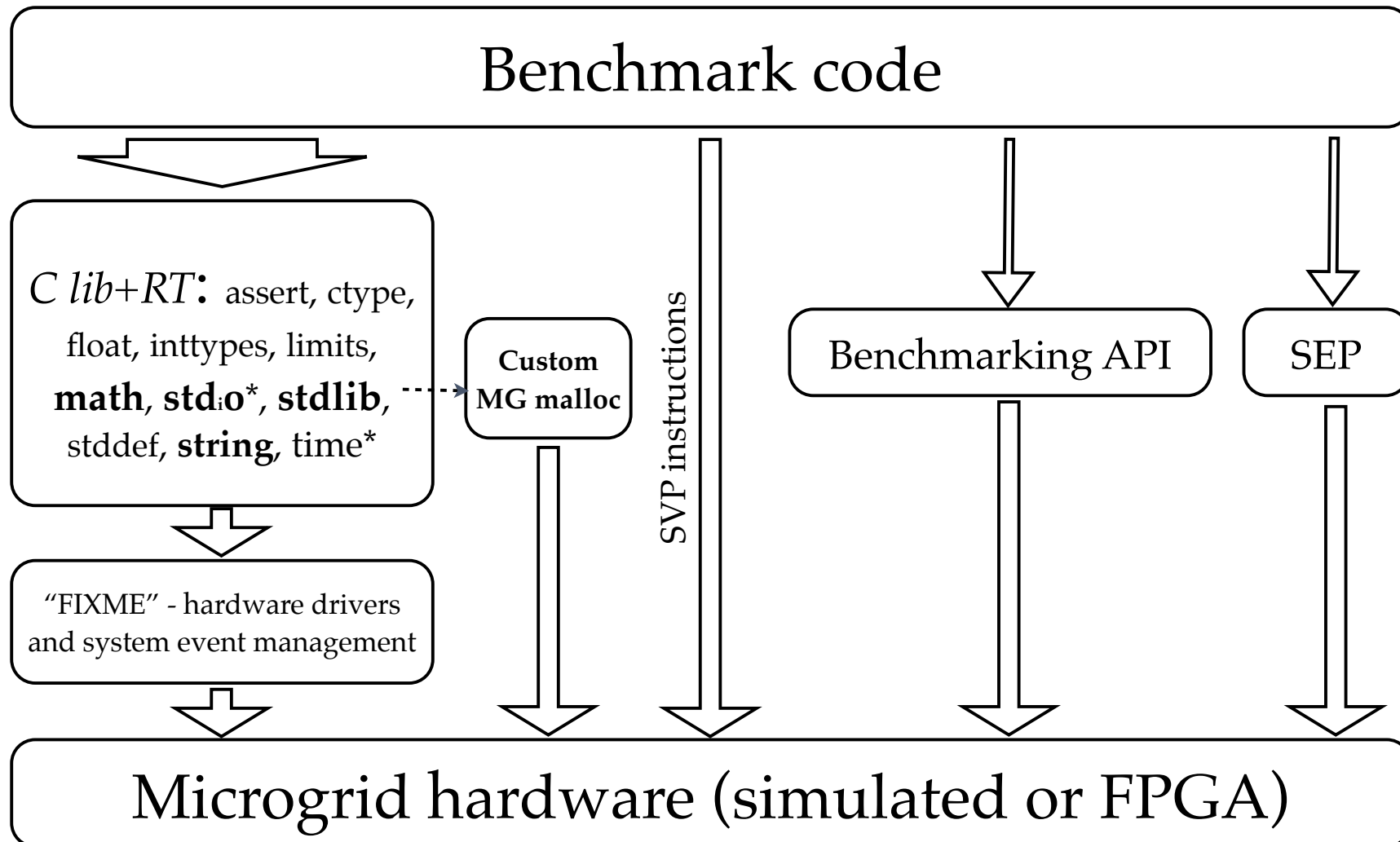
MGSim

UTLEON3

Distributed memory  
hw multithreaded  
multi-cores

Microgrid hardware model

# MGOS: SOFTWARE INTERFACES





# INSIDE THE MGOS

Component	Objects	Functions	Code words
Libc-dtoa-strtod	44	74	14528
Libc-math	165	238	13424
Libc-stdout	13	13	5040
Libc-dlmalloc	2	19	3600
Libc-string	26	26	1968
Libc-misc	9	28	1728
mginit	2	2	1120
Libmg-benchmarks	1	7	1072
Libmg-SEP	1	4	848
Libc-malloc-mg	2	9	832
Libsl-misc	3	12	736

Many copyright holders...

... but only one import tree: FreeBSD

And comparatively little MG-specific code!

