

APPLE-CORE: SVP AND MICROGRIDS

CAN WE STILL RETHINK THE HARDWARE/SOFTWARE INTERFACE IN GENERAL-PURPOSE PROCESSORS?

RAPHAEL 'KENA' POSS
UNIVERSITY OF AMSTERDAM, THE NETHERLANDS

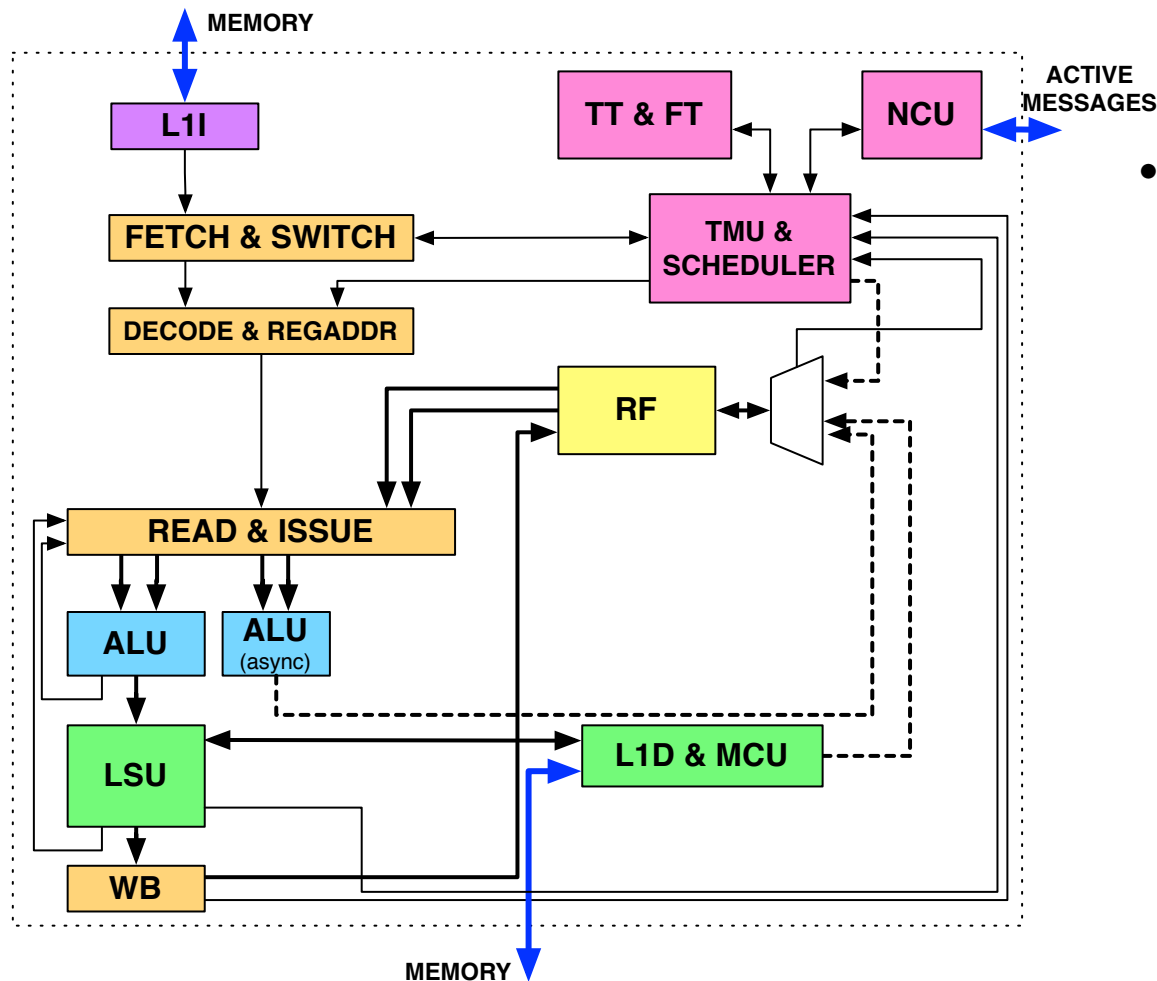
DSD 2012
CESME, IZMIR, TURKEY
SEPTEMBER 6TH, 2012



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?

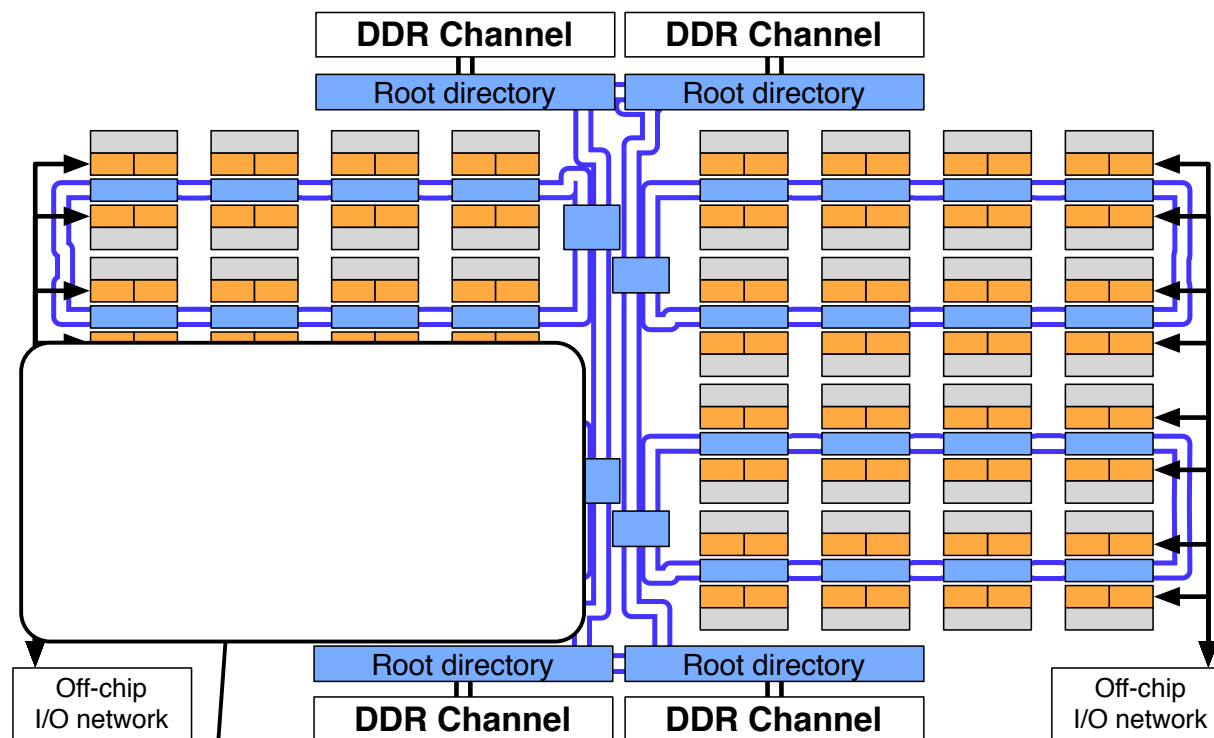
MICROGRIDS OF D-RISC CORES



- 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

EXAMPLE 128-CORE MICROGRID



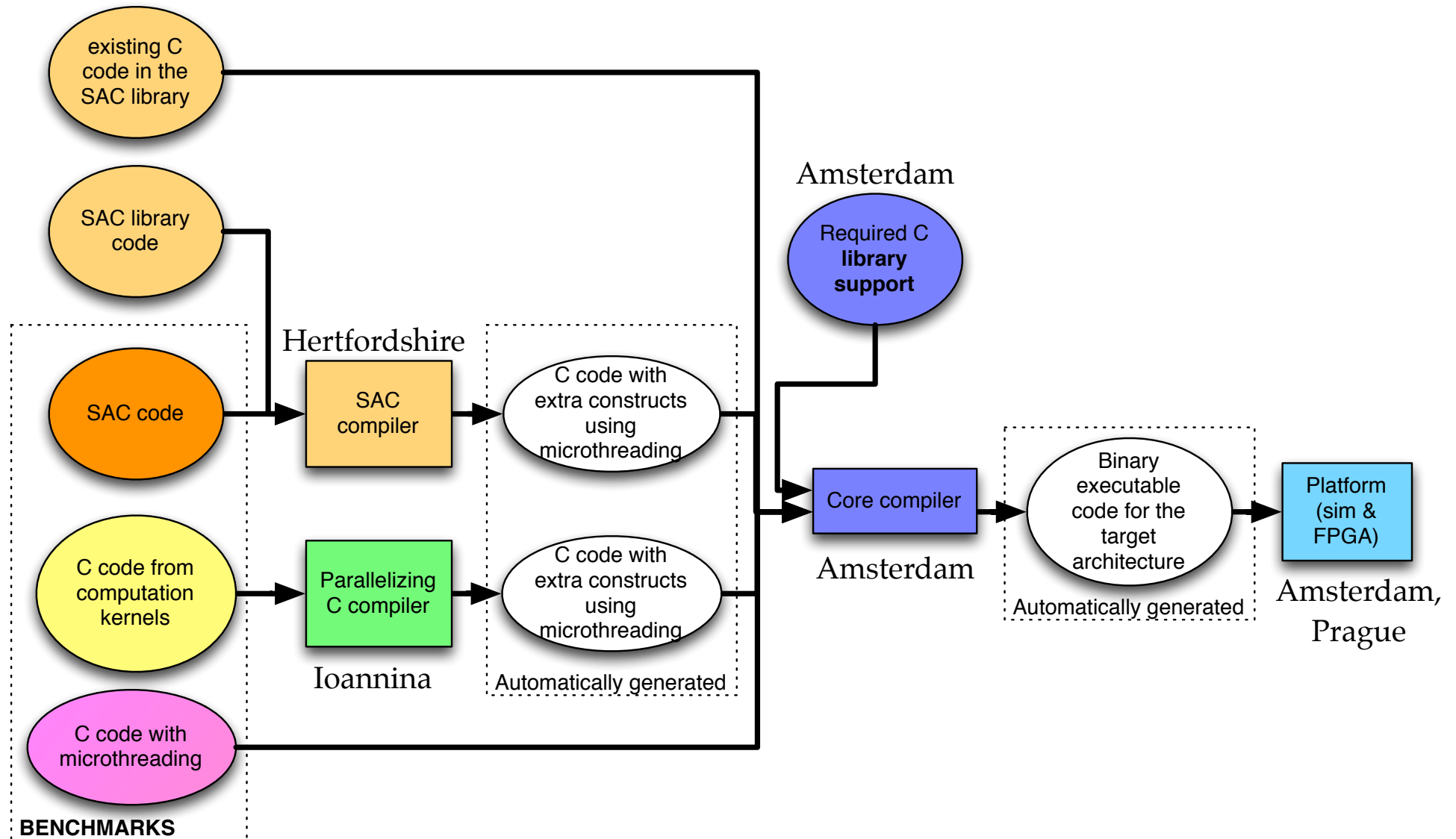
- 32000+ hw threads
- 5MB distributed cache
- shared MMU
= single virtual address space, protection using capabilities
- Weak cache coherency
- no support for global memory atomics – instead synchronization using point-to-point messaging

Area estimates with CACTI: 100mm² @ 35nm

A PERSPECTIVE SHIFT

CORE I7	Function call with 4 registers spilled 30-100 cycles	Predictable loop requires branch predictor + cache prefetching to maximize utilization 1+ cycles / iteration overhead
D-RISC WITH TMU IN HARDWARE	Bulk thread creation of 1 thread, 31 “fresh” registers ~15 cycles (7c sync, ~8c async)	Thread family 1 thread / “iteration” reuses common TMU and pipeline no BP nor prefetch needed 0+ iteration overhead

THE APPLE-CORE SOFTWARE STRATEGY



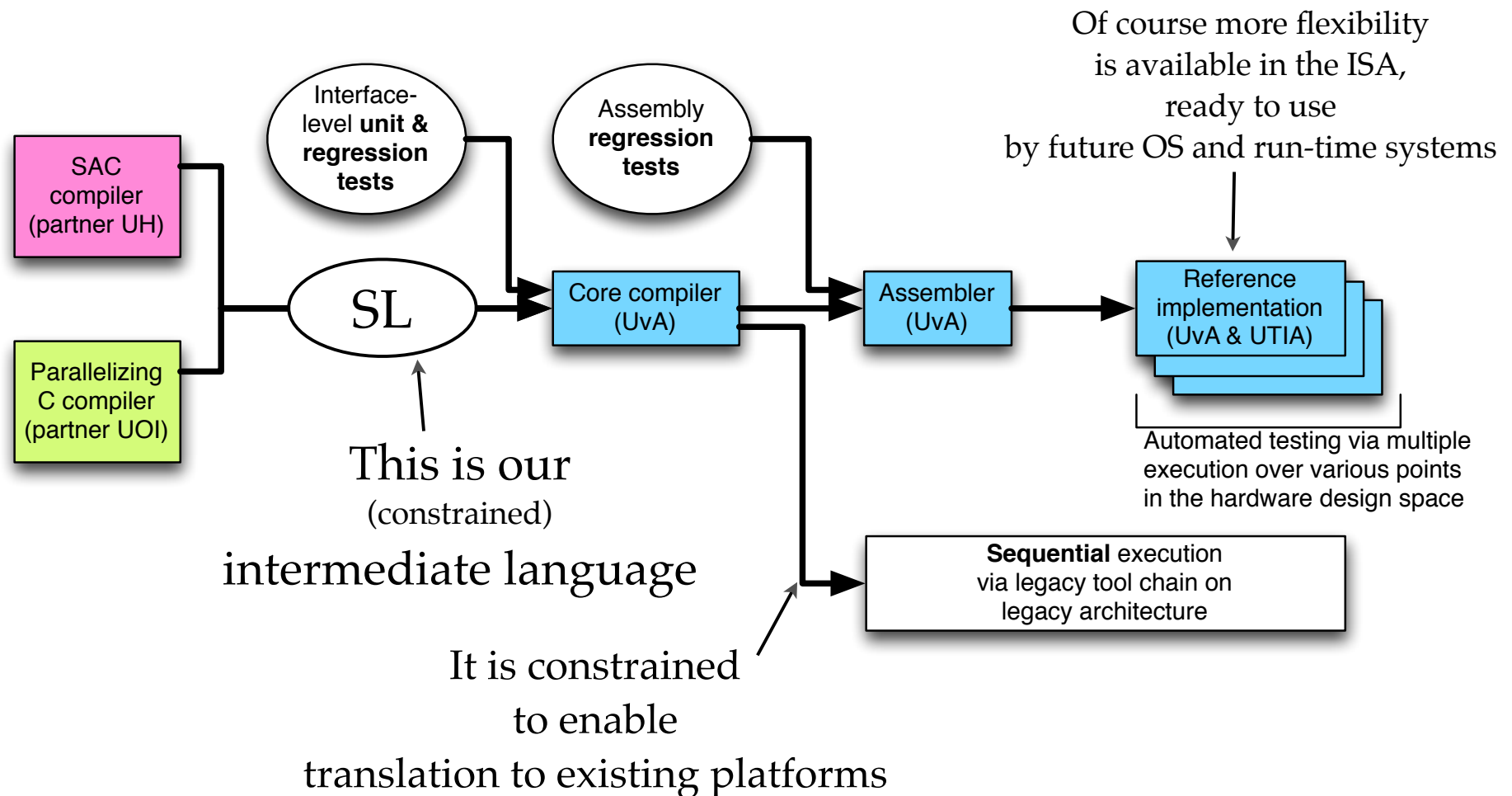
THE “MAIN” ISSUES UNCOVERED IN APPLE-CORE

- **Validation:** how to 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 exacerbated 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



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

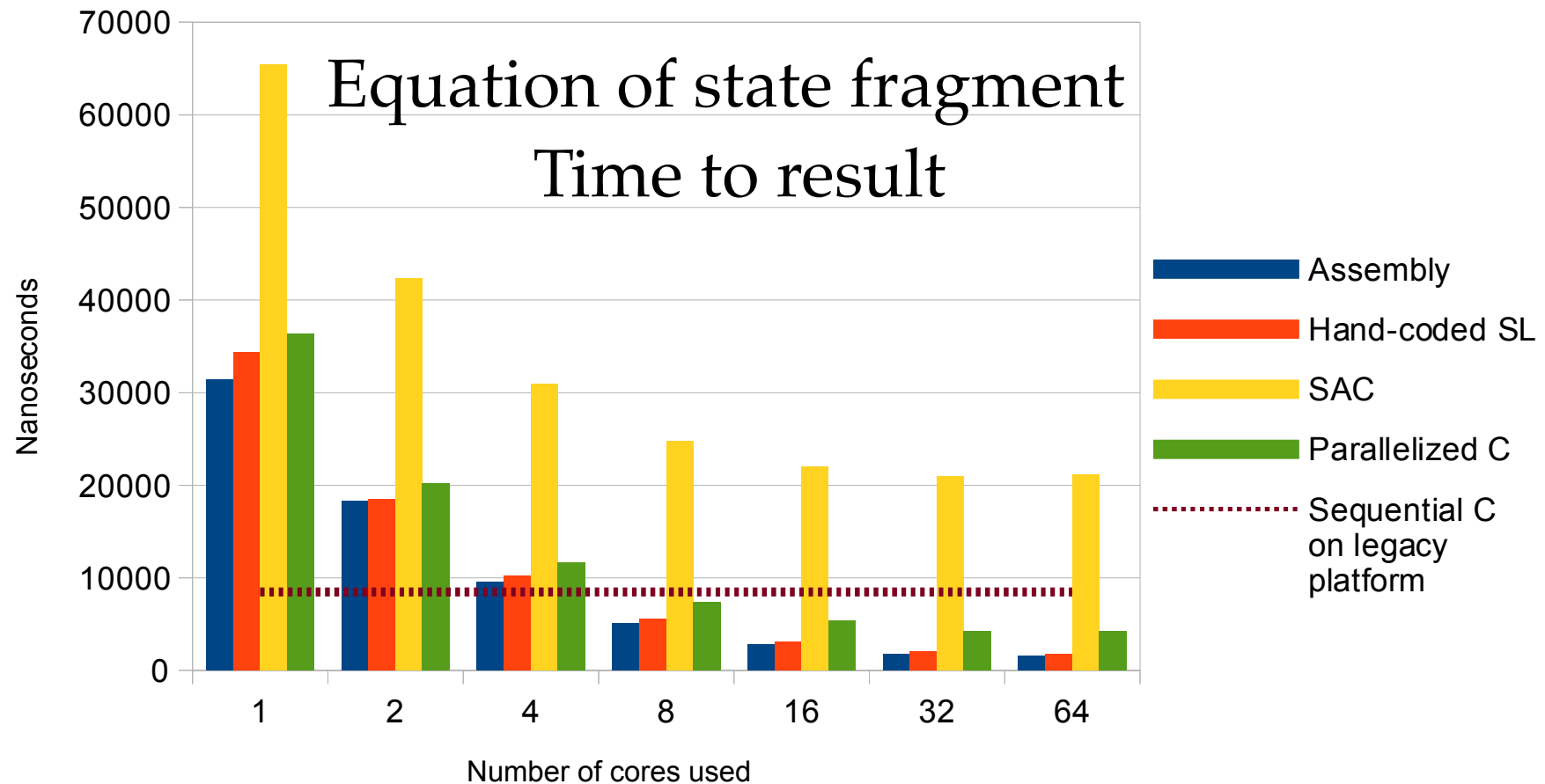
RESOURCE MANAGEMENT

- **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

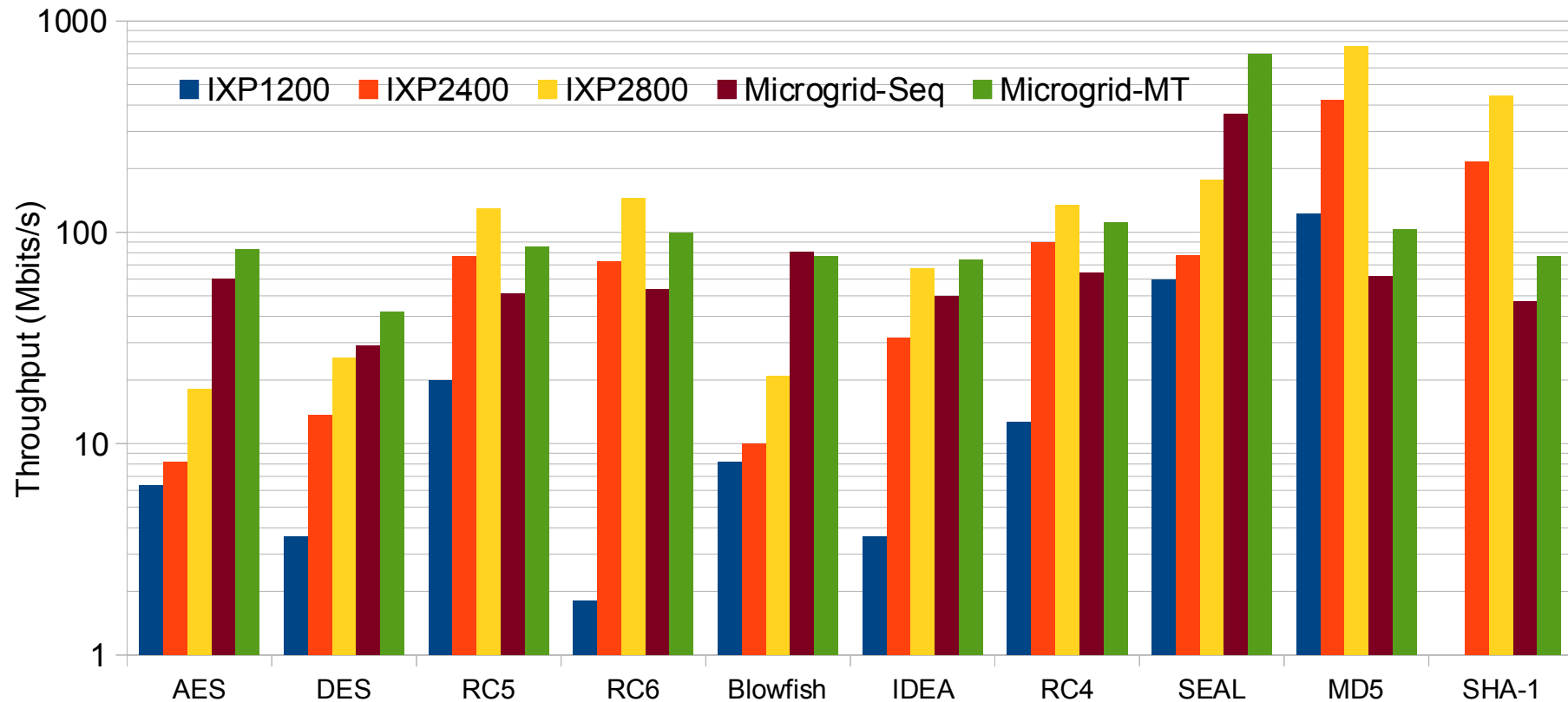


Legacy platform = MacBook Pro, Core 2 Duo @ 2.4GHz

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?

- ✓ **built enough infrastructure to fit the F/OSS landscape**
 - yet can't reuse most existing OS code: *no interrupts, no traps*
- ✓ **as planned, higher performance per area and per watt**
 - via hand-coded benchmarks: *granularity in SPEC is too coarse*
- **Follow-up research areas:**
 - *Internal* issues: memory consistency, scalable cache protocols, ISA semantics, etc.
 - *External* issues from outside architecture: how to virtualize? how to place tasks over so many “workers”? how to port existing OS code?
 - *Fundamental* issues: concurrent complexity theory?

THANK YOU!

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

SVP 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

EXTRA -

A PERSPECTIVE SHIFT

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