

The Scalable Commutativity Rule: Designing Scalable Software for Multicore Processors

Phil Gibbons

15-712 F15

Lecture 18

Today's Reminders

- **Project proposals due Friday midnight**
 - Email Kevin and me
- **No Class Monday**
- **Midterm on Wednesday**
 - Will cover assigned papers for 9/11 through 10/21
 - Understand high-level concepts & compare ideas/approaches across papers

2

The Scalable Commutativity Rule: Designing Scalable Software for Multicore Processors

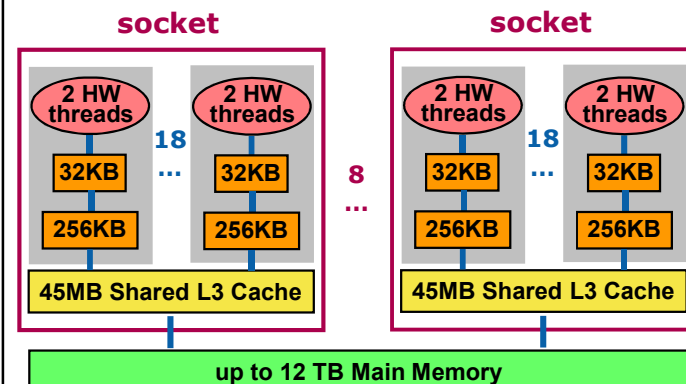
[SOSP'13 best paper]

- **Austin Clements** (Google)
- **Frans Kaashoek** (MIT)
- **Nickolai Zeldovich** (MIT)
- **Robert Morris** (MIT)
- **Eddie Kohler** (Harvard)



3

Multicore: 144-core Xeon Haswell E7-v3



Attach: Hard Drives & Flash Devices

4

An Analysis of Linux Scalability to Many Cores [OSDI'10]

- Analyzed 7 system apps running on Linux on 48-cores
 - Exim, memcached, Apache, PostgreSQL, gmake, Psearchy, MapReduce
 - Used an in-memory file system to explore non-disk limitations
- Identified & sought to remove all scalability bottlenecks
- Kernel changes are all localized
 - typically involve avoiding locks & atomic instructions by organizing data structures to avoid unnecessary sharing

5

Serializing Actions

Scalability bottlenecks can arise when tasks:

- Lock a shared data structure
- Write a shared memory location
- Compete for on-chip cache space
- Compete for on-chip interconnect or DRAM interface
- Are already mostly idle

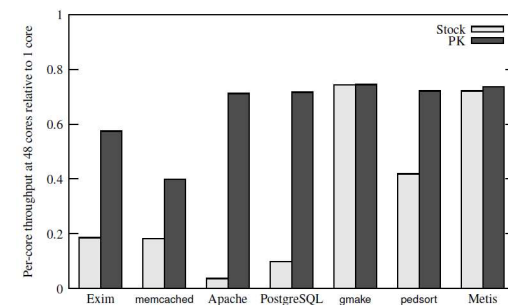
6

Scalability Problems & Fixes

Parallel accept		Apache
Concurrent accept system calls contend on shared socket fields.	⇒ User per-core backlog queues for listening sockets.	
entry reference counting		Apache, Exim
File name resolution contends on directory entry reference counts.	⇒ Use sloppy counters to reference count directory entry objects.	
Mount point (vfs mount) reference counting		Apache, Exim
Walking file name paths contends on mount point reference counts.	⇒ Use sloppy counters for mount point objects.	
IP packet destination (dst_entry) reference counting		memcached, Apache
IP packet transmission contends on routing table entries.	⇒ Use sloppy counters for IP routing table entries.	
Protocol memory usage tracking		memcached, Apache
Cores contend on counters for tracking protocol memory consumption.	⇒ Use sloppy counters for protocol usage counting.	
Acquiring directory entry (dentry) spin locks		Apache, Exim
Walking file name paths contends on per-directory entry spin locks.	⇒ Use a lock-free protocol in dlookup for checking filename matches.	
Mount point table spin lock		Apache, Exim
Resolving path names to mount points contends on a global spin lock.	⇒ Use per-core mount table caches.	
Adding files to the open list		Apache, Exim
Cores contend on a per-super block list that tracks open files.	⇒ Use per-core open file lists for each super block that has open files.	
Allocating DMA buffers		memcached, Apache
DMA memory allocations contend on the memory node 0 spin lock.	⇒ Allocate Ethernet device DMA buffers from the local memory node.	
False sharing in net_device and device		memcached, Apache, PostgreSQL
False sharing causes contention for read-only structure fields.	⇒ Place read-only fields on their own cache lines.	
False sharing in page		Exim
False sharing causes contention for read-mostly structure fields.	⇒ Place read-only fields on their own cache lines.	
inode lists		memcached, Apache
Cores contend on global locks protecting lists used to track inodes.	⇒ Avoid acquiring the locks when not necessary.	
Dcache lists		memcached, Apache
Cores contend on global locks protecting lists used to track dentries.	⇒ Avoid acquiring the locks when not necessary.	
Per-inode mutex		PostgreSQL
Cores contend on a per-inode mutex in lseek.	⇒ Use atomic reads to eliminate the need to acquire the mutex.	
Super-page fine grained locking		Metis
Super-page soft page faults contend on a per-process mutex.	⇒ Protect each super-page memory mapping with its own mutex.	
Zeroing super-pages		Metis
Zeroing super-pages flushes the contents of on-chip caches.	⇒ Use non-caching instructions to zero the contents of super-pages.	

7

Throughput: Before & After



Remaining Bottlenecks are not due to Linux

Application	Bottleneck
Exim	App: Contention on spool directories
memcached	HW: Transmit queues on NIC
Apache	HW: Receive queues on NIC
PostgreSQL	App: Application-level spin lock
gmake	App: Serial stages and stragglers
pedsort	HW: Cache capacity
Metis	HW: DRAM throughput

8

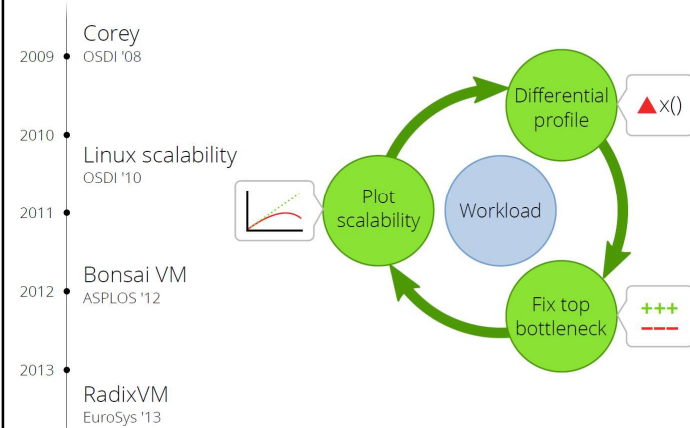
[Slides from SOSP'13 talk]

The Scalable Commutativity Rule: Designing Scalable Software for Multicore Processors

Austin T. Clements
M. Frans Kaashoek
Nikolai Zeldovich
Robert Morris
Eddie Kohler †

MIT CSAIL and † Harvard

Current approach to scalable software development



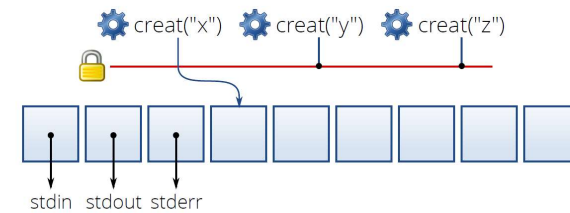
Current approach to scalable software development

Successful in practice because it focuses developer effort

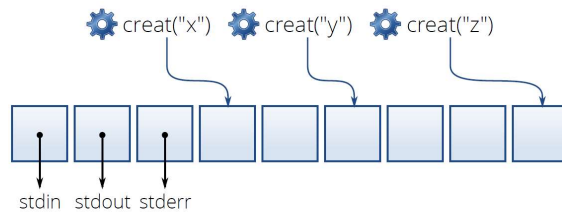
Disadvantages

- New workloads expose new bottlenecks
- More cores expose new bottlenecks
- The real bottlenecks may be in the interface design

Interface scalability example



Interface scalability example



Approach: Interface-driven scalability

The scalable commutativity rule

Whenever interface operations commute, they can be implemented in a way that scales.

Approach: Interface-driven scalability

The scalable commutativity rule

Whenever interface operations commute, they can be implemented in a way that scales.

creat with lowest FD

Commutates	Scalable implementation exists
?	
creat → 3	
creat → 4	

Approach: Interface-driven scalability

The scalable commutativity rule

Whenever interface operations commute, they can be implemented in a way that scales.

creat with lowest FD

Commutates	Scalable implementation exists
×	
?	
creat → 42	
creat → 17	

Approach: Interface-driven scalability

The scalable commutativity rule

**Whenever interface operations commute,
they can be implemented in a way that scales.**



Advantages of interface-driven scalability

The rule enables reasoning about scalability throughout the software design process

Design	Guides design of scalable interfaces
Implement	Sets a clear implementation target
Test	Systematic, workload-independent scalability testing

Contributions

The scalable commutativity rule

- Formalization of the rule and proof of its correctness
- State-dependent, interface-based commutativity

Commuter: An automated scalability testing tool

sv6: A scalable POSIX-like kernel

Outline

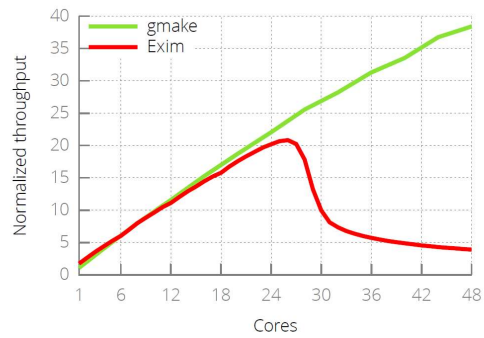
Defining the rule

- Definition of scalability
- Intuition
- Formalization

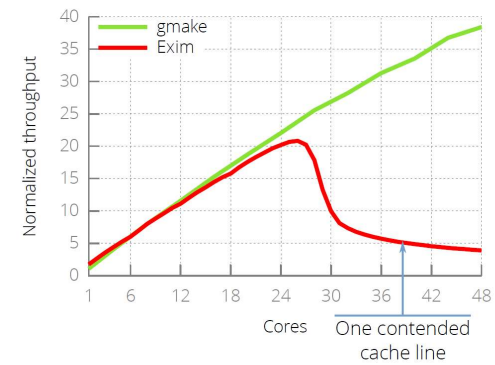
Applying the rule

- Commuter
- Evaluation

A scalability bottleneck

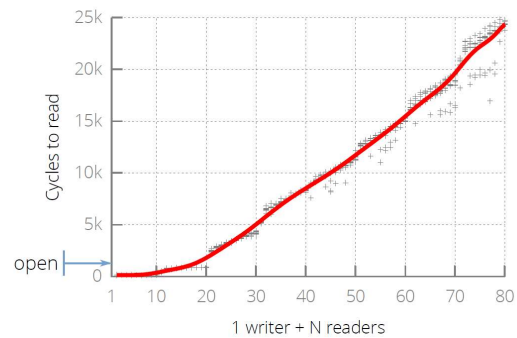


A scalability bottleneck



A single contended cache line can wreck scalability

Cost of a contended cache line



What scales on today's multicores?

		Core X		
		W	R	-
Core Y	W	✗	✗	✓
	R	✗	✓	✓
	-	✓	✓	-

We say two or more operations are *scalable* if they are *conflict-free*.

The intuition behind the rule

Whenever interface operations commute,
they can be implemented in a way that scales.

- Operations commute
- ⇒ results independent of order
- ⇒ communication is unnecessary
- ⇒ without communication, no conflicts

Formalizing the rule

Y SI-commutes in $X \parallel Y$ $\hat{=}$
 $\forall Y' \in \text{reorderings}(Y), Z: X \parallel Y \parallel Z \in \mathcal{P} \Leftrightarrow X \parallel Y' \parallel Z \in \mathcal{P}$

Y SIM-commutes in $X \parallel Y$ $\hat{=}$
 $\forall P \in \text{prefixes}(\text{reorderings}(Y)): P$ SI-commutes in $X \parallel P$.

An implementation m is a step function: $\text{state} \times \text{inv} \mapsto \text{state} \times \text{resp}$.

Given a specification \mathcal{P} ,
a history $X \parallel Y$ in which Y SIM-commutes,
and a reference implementation M that can generate $X \parallel Y$,
 \exists an implementation m of \mathcal{P} whose steps in Y are conflict-free.

Proof by simulation construction.

Formalizing the rule

Y SI-commutes in $X \parallel Y$ $\hat{=}$
 $\forall Y' \in \text{reorderings}(Y), Z: X \parallel Y \parallel Z \in \mathcal{P} \Leftrightarrow X \parallel Y' \parallel Z \in \mathcal{P}$

Y SIM-commutes in $X \parallel Y$ $\hat{=}$
 $\forall P \in \text{prefixes}(\text{reorderings}(Y)): P$ SI-commutes in $X \parallel P$.

An implementation m is a step function: $\text{state} \times \text{inv} \mapsto \text{state} \times \text{resp}$.

Given a specification \mathcal{P} ,
a history $X \parallel Y$ in which Y SIM-commutes,
and a reference implementation M that can generate $X \parallel Y$,
 \exists an implementation m of \mathcal{P} whose steps in Y are conflict-free.

Proof by simulation construction.

Commutativity is sensitive to
operations, arguments, and state

Example of using the rule

	Commutates	Scalable implementation exists
P1: creat P1: creat	✗	
P1: creat("/tmp/x") P2: creat("/etc/y")	✓	✓ (Linux)
P1: creat("/x") P2: creat("/y")	✓	✓
P1: creat("x", O_EXCL) P2: creat("x", O_EXCL)		
Same CWD	✗	
Different CWD	✓	✓

Input: Symbolic model

```

SymInode   = tstruct(data = tlist(SymByte),
                    nlink = SymInt)
SymIMap    = tdict(SymInt, SymInode)
SymFilename = tuninterpreted('Filename')
SymDir     = tdict(SymFilename, SymInt)

class POSIX:
    def __init__(self):
        self.fname_to_inum = SymDir.any()
        self.inodes = SymIMap.any()

    @symargs(src=SymFilename, dst=SymFilename)
    def rename(self, src, dst):
        if src not in self.fname_to_inum:
            return (-1, errno.ENOENT)
        if src == dst:
            return 0
        if dst in self.fname_to_inum:
            self.inodes[self.fname_to_inum[dst]].nlink -= 1
        self.fname_to_inum[dst] = self.fname_to_inum[src]
        del self.fname_to_inum[src]
        return 0

```

Symbolic model



Commutativity conditions

```

@symargs(src=SymFilename, dst=SymFilename)
def rename(self, src, dst):
    if src not in self.fname_to_inum:
        return (-1, errno.ENOENT)
    if src == dst:
        return 0
    if dst in self.fname_to_inum:
        self.inodes[self.fname_to_inum[dst]].nlink -= 1
    self.fname_to_inum[dst] = self.fname_to_inum[src]
    del self.fname_to_inum[src]
    return 0

```

rename(a, b) and rename(c, d) commute if:

- Both source files exist and all names are different
- Neither source file exists
- a xor c exists, and it is not the other rename's destination
- Both calls are self-renames
- One call is a self-rename of an existing file and a != c
- a & c are hard links to the same inode, a != c, and b == d

Symbolic model



Test cases

rename(a, b) and rename(c, d) commute if:

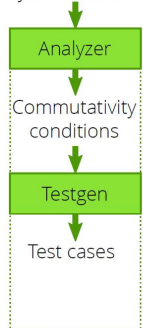
- Both source files exist and all names are different
- Neither source file exists
- a xor c exists, and it is not the other rename's destination
- Both calls are self-renames
- One call is a self-rename of an existing file and a != c
- a & c are hard links to the same inode, a != c, and b == d

```

void setup() {
    close(creat("f0", 0666));
    close(creat("f2", 0666));
}
void test_opA() { rename("f0", "f1"); }
void test_opB() { rename("f2", "f3"); }

```

Symbolic model

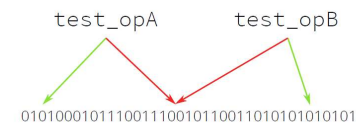


Output: Conflicting cache lines

```

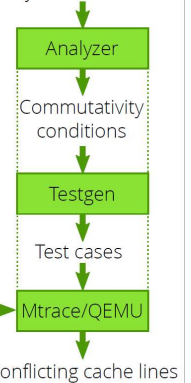
void setup() {
    close(creat("f0", 0666));
    close(creat("f2", 0666));
}
void test_opA() { rename("f0", "f1"); }
void test_opB() { rename("f2", "f3"); }

```



Linux

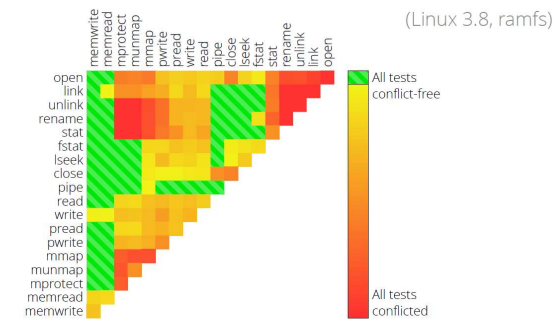
Symbolic model



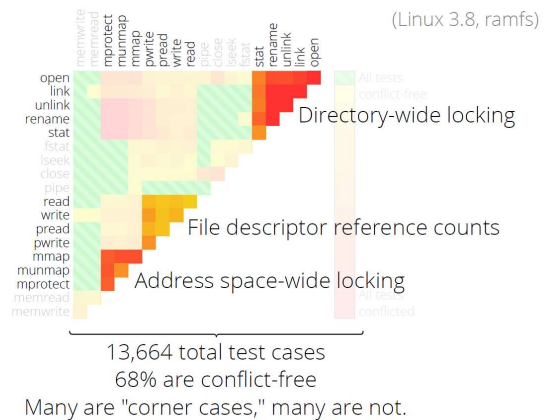
Evaluation

Does the rule help build scalable systems?

Commuter finds non-scalable cases in Linux



Commuter finds non-scalable cases in Linux



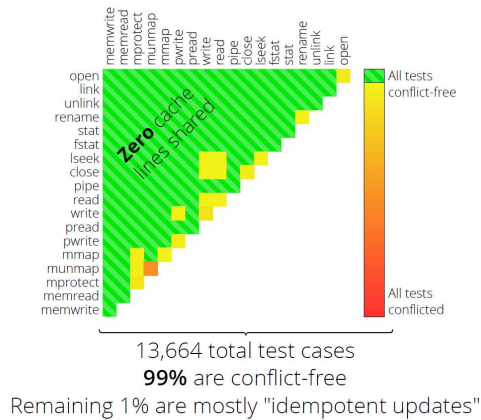
sv6: A scalable OS

POSIX-like operating system

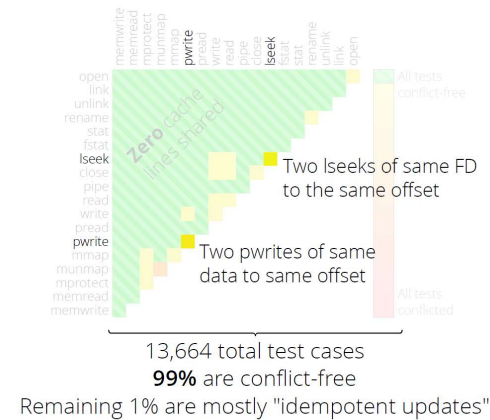
File system and virtual memory system follow commutativity rule

Implementation using standard parallel programming techniques,
but guided by Commuter

Commutative operations can be made to scale



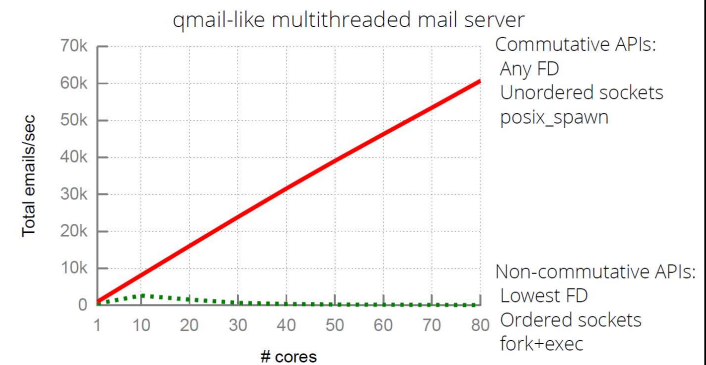
Commutative operations can be made to scale



Refining POSIX with the rule

- Lowest FD versus any FD
- stat versus xstat
- Unordered sockets
- Delayed munmap
- fork+exec versus posix_spawn

Commutative operations matter to app scalability



Related work

Commutativity and concurrency

- [Bernstein '81]
- [Weihl '88]
- [Steele '90]
- [Rinard '97]
- [Shapiro '11]

Laws of Order [Attiya '11]

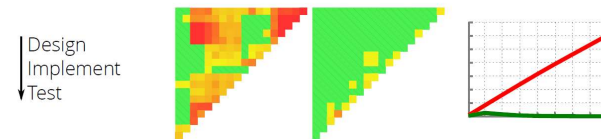
Disjoint-access parallelism [Israeli '94]

Scalable locks [MCS '91]

Scalable reference counting [Ellen '07, Corbet '10]

Conclusion

Whenever interface operations commute,
they can be implemented in a way that scales.



Check it out at <http://pdos.csail.mit.edu/commuter>

[End of slides from SOSP'13 talk]

Discussion

- Scalable Commutativity Rule only implies that **there exists an implementation with conflict-free accesses**
 - Implementation constructed in proof is not practical

- Can scalability suffer even for conflict-free accesses?

1. Lock a shared data structure
2. Write a shared memory location
3. Compete for on-chip cache space
4. Compete for on-chip interconnect or DRAM interface
5. Are already mostly idle

Yes. Only addresses first two sources of lack of scalability

43

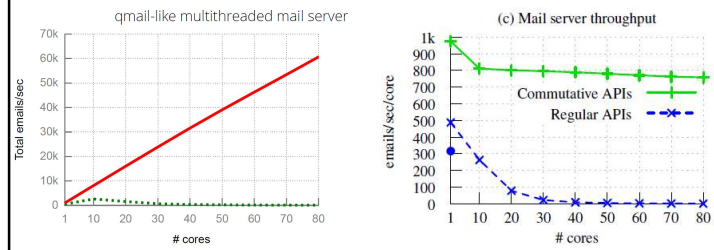
Interface Changes for POSIX

1. Decompose compound operations
2. Embrace specification non-determinism
3. Permit weak ordering
4. Release resources asynchronously

- Lowest FD versus any FD (2)
- stat versus xstat (1)
- Unordered sockets (3)
- Delayed munmap (4)
- fork+exec versus posix_spawn (1)

44

Application Performance



45

Next Class

Wednesday's Midterm

46