THE OS SCHEDULER: A PERFORMANCE-CRITICAL COMPONENT IN LINUX CLUSTER ENVIRONMENTS

KEYNOTE FOR BPOE-9 *OASPLOS2018* THE NINTH WORKSHOP ON BIG DATA BENCHMARKS, PERFORMANCE, OPTIMIZATION AND EMERGING HARDWARE

By Jean-Pierre Lozi Oracle Labs

CLUSTER COMPUTING

Multicore servers with dozens of cores

- Common for e.g., a hadoop cluster, a distributed graph analytics engine, multiple apps...
- High cost of infrastructure, high energy consumption



CLUSTER COMPUTING

Multicore servers with dozens of cores

- Common for e.g., a hadoop cluster, a distributed graph analytics engine, multiple apps...
- High cost of infrastructure, high energy consumption

Linux-based software stack

Low (license) cost, yet high reliability



CLUSTER COMPUTING

Multicore servers with dozens of cores

- Common for e.g., a hadoop cluster, a distributed graph analytics engine, multiple apps...
- High cost of infrastructure, high energy consumption

Linux-based software stack

- Low (license) cost, yet high reliability
- Challenge: don't waste cycles!
- Reduces infrastructure and energy costs
- Improves bandwidth and latency





3







reducing network usage (e.g. HDFS): common optimizations 3





It must be! 15 years ago, Linus Torvalds was already saying:

"And you have to realize that there are not very many things that have aged as well as the scheduler. Which is just another proof that scheduling is easy."

It must be! 15 years ago, Linus Torvalds was already saying:

"And you have to realize that there are not very many things that have aged as well as the scheduler. Which is just another proof that scheduling is easy."

Since then, people have been running applications on their multicore machines all the time, and they run, CPU usage is high, everything seems fine.

It must be! 15 years ago, Linus Torvalds was already saying:

"And you have to realize that there are not very many things that have aged as well as the scheduler. Which is just another proof that scheduling is easy."

Since then, people have been running applications on their multicore machines all the time, and they run, CPU usage is high, everything seems fine.

But would you notice if some cores remained idle intermittently, when they shouldn't?

It must be! 15 years ago, Linus Torvalds was already saying:

"And you have to realize that there are not very many things that have aged as well as the scheduler. Which is just another proof that scheduling is easy."

Since then, people have been running applications on their multicore machines all the time, and they run, CPU usage is high, everything seems fine.

- But would you notice if some cores remained idle intermittently, when they shouldn't?
- Do you keep monitoring tools (htop) running all the time?

It must be! 15 years ago, Linus Torvalds was already saying:

"And you have to realize that there are not very many things that have aged as well as the scheduler. Which is just another proof that scheduling is easy."

Since then, people have been running applications on their multicore machines all the time, and they run, CPU usage is high, everything seems fine.

- But would you notice if some cores remained idle intermittently, when they shouldn't?
- Do you keep monitoring tools (htop) running all the time?
- Even if you do, would you be able to identify faulty behavior from normal noise?

It must be! 15 years ago, Linus Torvalds was already saying:

"And you have to realize that there are not very many things that have aged as well as the scheduler. Which is just another proof that scheduling is easy."

Since then, people have been running applications on their multicore machines all the time, and they run, CPU usage is high, everything seems fine.

- But would you notice if some cores remained idle intermittently, when they shouldn't?
- Do you keep monitoring tools (htop) running all the time?
- Even if you do, would you be able to identify faulty behavior from normal noise?
- Would you ever suspect the scheduler?

 Over the past few years of working on various projects, we sometimes saw strange, hard to explain performance results.

 Over the past few years of working on various projects, we sometimes saw strange, hard to explain performance results.

• An example: running a TPC-H benchmark on a 64-core machine, our runs much faster when pinning threads to cores than when we let the Linux scheduler do its job.

 Over the past few years of working on various projects, we sometimes saw strange, hard to explain performance results.

- An example: running a TPC-H benchmark on a 64-core machine, our runs much faster when pinning threads to cores than when we let the Linux scheduler do its job.
- Memory locality issue? Impossible, hardware counters showed no difference in the % of remote memory accesses, in cache misses, etc.

 Over the past few years of working on various projects, we sometimes saw strange, hard to explain performance results.

- An example: running a TPC-H benchmark on a 64-core machine, our runs much faster when pinning threads to cores than when we let the Linux scheduler do its job.
 - Memory locality issue? Impossible, hardware counters showed no difference in the % of remote memory accesses, in cache misses, etc.
 - Contention over some resource (spinlock, etc.)? We investigated this for a long time, but couldn't find anything that looked off.

 Over the past few years of working on various projects, we sometimes saw strange, hard to explain performance results.

- An example: running a TPC-H benchmark on a 64-core machine, our runs much faster when pinning threads to cores than when we let the Linux scheduler do its job.
 - Memory locality issue? Impossible, hardware counters showed no difference in the % of remote memory accesses, in cache misses, etc.
 - Contention over some resource (spinlock, etc.)? We investigated this for a long time, but couldn't find anything that looked off.
 - Overhead of context switches? Threads moved a lot but we proved that the overhead was negligible.

 Over the past few years of working on various projects, we sometimes saw strange, hard to explain performance results.

- An example: running a TPC-H benchmark on a 64-core machine, our runs much faster when pinning threads to cores than when we let the Linux scheduler do its job.
 - Memory locality issue? Impossible, hardware counters showed no difference in the % of remote memory accesses, in cache misses, etc.
 - Contention over some resource (spinlock, etc.)? We investigated this for a long time, but couldn't find anything that looked off.
 - Overhead of context switches? Threads moved a lot but we proved that the overhead was negligible.
- We ended up suspecting the core behavior of the scheduler.

 Over the past few years of working on various projects, we sometimes saw strange, hard to explain performance results.

- An example: running a TPC-H benchmark on a 64-core machine, our runs much faster when pinning threads to cores than when we let the Linux scheduler do its job.
 - Memory locality issue? Impossible, hardware counters showed no difference in the % of remote memory accesses, in cache misses, etc.
 - Contention over some resource (spinlock, etc.)? We investigated this for a long time, but couldn't find anything that looked off.
 - Overhead of context switches? Threads moved a lot but we proved that the overhead was negligible.
- We ended up suspecting the core behavior of the scheduler.
 - We implemented high-resolution tracing tools and saw that some cores were idle while others overloaded...

 Over the past few years of working on various projects, we sometimes saw strange, hard to explain performance results.



Slowed down execution

• This is how we found our first performance bug. Which made us investigate more...

- This is how we found our first performance bug. Which made us investigate more...
- In the end: four Linux scheduler performance bugs that we found and analyzed

- This is how we found our first performance bug. Which made us investigate more...
- In the end: four Linux scheduler performance bugs that we found and analyzed
- Always the same symptom: idle cores while others are overloaded
- The bug-hunting was tough, and led us to develop our own tools

- This is how we found our first performance bug. Which made us investigate more...
- In the end: four Linux scheduler performance bugs that we found and analyzed
- Always the same symptom: idle cores while others are overloaded
 - The bug-hunting was tough, and led us to develop our own tools
- Performance overhead of some of the bugs :
 - 12-23% performance improvement on a popular database with TPC-H
 - 137× performance improvement on HPC workloads

- This is how we found our first performance bug. Which made us investigate more...
- In the end: four Linux scheduler performance bugs that we found and analyzed
- Always the same symptom: idle cores while others are overloaded
 - The bug-hunting was tough, and led us to develop our own tools
- Performance overhead of some of the bugs :
 - 12-23% performance improvement on a popular database with TPC-H
 - 137× performance improvement on HPC workloads
- Not always possible to provide a simple, working fix...
- Intrisic problems with the design of the scheduler?

Main takeaway of our analysis: more research must be directed towards implementing an efficient scheduler for multicore architectures, because contrary to what a lot of us think, this is *not* a solved problem!

Main takeaway of our analysis: more research must be directed towards implementing an efficient scheduler for multicore architectures, because contrary to what a lot of us think, this is *not* a solved problem!

Need convincing? Let's go through it together...

Main takeaway of our analysis: more research must be directed towards implementing an efficient scheduler for multicore architectures, because contrary to what a lot of us think, this is *not* a solved problem!

Need convincing? Let's go through it together...

...starting with a bit of background...

The Linux Scheduler: a Decade of Wasted Cores

Jean-Pierre Lozi Université Nice Sophia-Antipolis jplozi@unice.fr

Fabien Gaud

Coho Data

me@fabiengaud.net

Baptiste Lepers EPFL baptiste.lepers@epfl.ch

Vivien Quéma

Grenoble INP / ENSIMAG

vivien.guema@imag.fr

Justin Funston University of British Columbia jfunston@ece.ubc.ca

Alexandra Fedorova University of British Columbia sasha@ece.ubc.ca

Main takeaway of our analysis: more research must be directed towards implementing an efficient scheduler for multicore architectures, because contrary to what a lot of us think, this is *not* a solved problem!

Need convincing? Let's go through it together...

...starting with a bit of background...

THE OS SCHEDULER: A PERFORMANCE-CRITICAL COMPONENT IN LINUX CLUSTER ENVIRONMENTS 7

THIS TALK

THE COMPLETELY FAIR SCHEDULER (CFS): CONCEPT



THE COMPLETELY FAIR SCHEDULER (CFS): CONCEPT

One runqueue where threads are globally sorted by *runtime*




One runqueue where threads are globally sorted by *runtime*

Core 0



One runqueue where threads are globally sorted by *runtime*

Core 0







• One runqueue per core to avoid contention



- One runqueue per core to avoid contention
- CFS periodically balances "loads":
- $load(task) = weight^1 x \% cpu use^2$

¹The lower the niceness, the higher the weight



- One runqueue per core to avoid contention
- CFS periodically balances "loads":

$load(task) = weight^1 x \% cpu use^2$

¹The lower the niceness, the higher the weight

²We don't want a high-priority thread that sleeps a lot to take a whole CPU for itself and then mostly sleep!



- One runqueue per core to avoid contention
- CFS periodically balances "loads":

$load(task) = weight^1 x \% cpu use^2$

¹The lower the niceness, the higher the weight

²We don't want a high-priority thread that sleeps a lot to take a whole CPU for itself and then mostly sleep!



Since there can be many cores: hierarchical approach!



















Note that only the average load of groups is considered

- Note that only the average load of groups is considered
- If for some reason the lower-level load-balancing fails, nothing happens at a higher level:

Note that only the average load of groups is considered



Note that only the average load of groups is considered

If for some reason the lower-level load-balancing fails, nothing happens at a higher level:
AVG(L)=3000
L=0
L=6000
L=3000
L=3000



Note that only the average load of groups is considered

• If for some reason the lower-level load-balancing fails, nothing happens at a higher level:



Note that only the average load of groups is considered

• If for some reason the lower-level load-balancing fails, nothing happens at a higher level:



Load calculations are actually more complicated, use more heuristics.

- Load calculations are actually more complicated, use more heuristics.
- One of them aims to increase fairness between "sessions".

- Load calculations are actually more complicated, use more heuristics.
- One of them aims to increase fairness between "sessions".
- **Objective:** making sure that launching lots of threads from one terminal doesn't prevent other processes on the machine (potentially from other users) from running.

- Load calculations are actually more complicated, use more heuristics.
- One of them aims to increase fairness between "sessions".
- **Objective:** making sure that launching lots of threads from one terminal doesn't prevent other processes on the machine (potentially from other users) from running.
- Otherwise, easy to use more resources than other users by spawning many threads...

- Load calculations are actually more complicated, use more heuristics.
- One of them aims to increase fairness between "sessions".





- Load calculations are actually more complicated, use more heuristics.
- One of them aims to increase fairness between "sessions".



- Load calculations are actually more complicated, use more heuristics.
- One of them aims to increase fairness between "sessions".



- Load calculations are actually more complicated, use more heuristics.
- One of them aims to increase fairness between "sessions".



- Load calculations are actually more complicated, use more heuristics.
- Solution: divide the load of a task by the number of threads in its tty...

- Load calculations are actually more complicated, use more heuristics.
- Solution: divide the load of a task by the number of threads in its tty...



- Load calculations are actually more complicated, use more heuristics.
- Solution: divide the load of a task by the number of threads in its tty...



- Load calculations are actually more complicated, use more heuristics.
- Solution: divide the load of a task by the number of threads in its tty...



- Load calculations are actually more complicated, use more heuristics.
- Solution: divide the load of a task by the number of threads in its tty...






L=1000

Load(thread) = %cpu × weight / #threads = 100 × 10 / 1 = 1000

Session (tty) 1

L=125	L=125
L=125	L=125
L=125	L=125
L=125	L=125

Session (tty) 2

Load(thread) = %cpu × weight / #threads = 100 × 10 / 8 = 125





















- Another example, on a 64-core machine, with load balancing:
 - First between pairs of cores (Bulldozer architecture, a bit like hyperthreading)
 - Then between NUMA nodes

Another example, on a 64-core machine, with load balancing:

- First between pairs of cores (Bulldozer architecture, a bit like hyperthreading)
- Then between NUMA nodes
- User 1 launches : ssh <machine> R & ssh <machine> R &

Another example, on a 64-core machine, with load balancing:

- First between pairs of cores (Bulldozer architecture, a bit like hyperthreading)
- Then between NUMA nodes
- User 1 launches :

ssh <machine> R & ssh <machine> R &

User 2 launches : ssh <machine> make -j 64 kernel

Another example, on a 64-core machine, with load balancing:

- First between pairs of cores (Bulldozer architecture, a bit like hyperthreading)
- Then between NUMA nodes
- User 1 launches : ssh <machine> R & ssh <machine> R &
- User 2 launches : ssh <machine> make -j 64 kernel
- The bug happens at two levels :
 - Other core on pair of core idle
- Other cores on NUMA node less busy...





- The bug happens at two levels :
 - Other core on pair of core idle
 - Other cores on NUMA node less busy...



1024



The bug happens at two levels :

0

- Other core on pair of core idle
- Other cores on NUMA node less busy...

• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average





• A simple solution: balance the *minimum load* of groups instead of the average





• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average









• A simple solution: balance the *minimum load* of groups instead of the average





• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average


• A simple solution: balance the *minimum load* of groups instead of the average



• A simple solution: balance the *minimum load* of groups instead of the average



After the fix, make runs 13% faster, and R is not impacted

• A simple solution: balance the *minimum load* of groups instead of the average



After the fix, make runs 13% faster, and R is not impacted

• A simple solution, but is it ideal? Minimum load more volatile than average...

• A simple solution: balance the *minimum load* of groups instead of the average



After the fix, make runs 13% faster, and R is not impacted

• A simple solution, but is it ideal? Minimum load more volatile than average...

May cause lots of unnecessary rebalancing. Revamping load calculations needed?

• Hierarchical load balancing is based on groups of cores named scheduling domains

• Hierarchical load balancing is based on groups of cores named scheduling domains

Based on affinity, i.e., pairs of cores, dies, CPUs, NUMA nodes, etc.

- Hierarchical load balancing is based on groups of cores named scheduling domains
 Based on affinity, i.e., pairs of cores, dies, CPUs, NUMA nodes, etc.
- Each scheduling domain contains groups that are the lower-level scheduling domains

- Hierarchical load balancing is based on groups of cores named scheduling domains
 Based on affinity, i.e., pairs of cores, dies, CPUs, NUMA nodes, etc.
- Each scheduling domain contains groups that are the lower-level scheduling domains
- For instance, on our 64-core AMD Bulldozer machine:

- Hierarchical load balancing is based on groups of cores named scheduling domains
 Based on affinity, i.e., pairs of cores, dies, CPUs, NUMA nodes, etc.
- Each scheduling domain contains groups that are the lower-level scheduling domains
- For instance, on our 64-core AMD Bulldozer machine:
- At level 1, each pair of core (scheduling domains) contain cores (scheduling groups)

- Hierarchical load balancing is based on groups of cores named scheduling domains
 Based on affinity, i.e., pairs of cores, dies, CPUs, NUMA nodes, etc.
- Each scheduling domain contains groups that are the lower-level scheduling domains
- For instance, on our 64-core AMD Bulldozer machine:
- At level 1, each pair of core (scheduling domains) contain cores (scheduling groups)
- At level 2, each CPU (s.d.) contain pairs of cores (s.g.)

- Hierarchical load balancing is based on groups of cores named scheduling domains
 Based on affinity, i.e., pairs of cores, dies, CPUs, NUMA nodes, etc.
- Each scheduling domain contains groups that are the lower-level scheduling domains
- For instance, on our 64-core AMD Bulldozer machine:
- At level 1, each pair of core (scheduling domains) contain cores (scheduling groups)
- At level 2, each CPU (s.d.) contain pairs of cores (s.g.)
- At level 3, each group of directly connected CPUs (s.d.) contain CPUs (s.g.)

- Hierarchical load balancing is based on groups of cores named scheduling domains
 Based on affinity, i.e., pairs of cores, dies, CPUs, NUMA nodes, etc.
- Each scheduling domain contains groups that are the lower-level scheduling domains
- For instance, on our 64-core AMD Bulldozer machine:
- At level 1, each pair of core (scheduling domains) contain cores (scheduling groups)
- At level 2, each CPU (s.d.) contain pairs of cores (s.g.)
- At level 3, each group of directly connected CPUs (s.d.) contain CPUs (s.g.)
- At level 4, the whole machine (s.d.) contains group of directly connected CPUs (s.g.)





At the **first level**, the first core balances load with the other core on the same pair (because they share resources, high affinity)





At the 3rd level, the first CPU balances load with directly connected CPUS



At the **4**th level, the first group of directly connected CPUs balances load with the other groups of directly connected CPUs





Groups of CPUs built by:

(2) picking first
CPU not in a
group and
looking for all
directly
connected CPUs



And then stop, because all CPUs are in a group





Suppose we taskset an application on these two CPUs, two hops apart (16 threads)



















• Fix: build the domains by creating one "directly connected" group for every CPU

Instead of the first CPU and the first one not "covered" by a group

- Fix: build the domains by creating one "directly connected" group for every CPU
 Instead of the first CPU and the first one not "covered" by a group
- Performance improvement of NAS applications on two nodes :

Application	With bug	After fix	Improvement
BT	99	56	1.75x
CG	42	15	2.73x
EP	73	36	2x
LU	1040	38	27x

- Fix: build the domains by creating one "directly connected" group for every CPU
 Instead of the first CPU and the first one not "covered" by a group
- Performance improvement of NAS applications on two nodes :

Application	With bug	After fix	Improvement
BT	99	56	1.75x
CG	42	15	2.73x
EP	73	36	2x
LU	1040	38	27x

- Very good improvement for LU because more threads than cores if can't use 16 cores
- Solves spinlock issues (incl. potential convoys)

BUG 3/4: MISSING SCHEDULING DOMAINS

In addition to this, when domains re-built, levels 3 and 4 not re-built...

BUG 3/4: MISSING SCHEDULING DOMAINS

In addition to this, when domains re-built, levels 3 and 4 not re-built...

I.e., no balancing between directly connected or 1-hop CPUs (i.e. any CPU)
In addition to this, when domains re-built, levels 3 and 4 not re-built...

- I.e., no balancing between directly connected or 1-hop CPUs (i.e. any CPU)
- Happens for instance when disabling and re-enabling a core

- In addition to this, when domains re-built, levels 3 and 4 not re-built...
 - I.e., no balancing between directly connected or 1-hop CPUs (i.e. any CPU)
 - Happens for instance when disabling and re-enabling a core
- Launch an application, first thread created on CPU 1

- In addition to this, when domains re-built, levels 3 and 4 not re-built...
 - I.e., no balancing between directly connected or 1-hop CPUs (i.e. any CPU)
 - Happens for instance when disabling and re-enabling a core
- Launch an application, first thread created on CPU 1
- First thread will stay on CPU 1, next threads will be created on CPU 1 (default Linux)

- In addition to this, when domains re-built, levels 3 and 4 not re-built...
 - I.e., no balancing between directly connected or 1-hop CPUs (i.e. any CPU)
 - Happens for instance when disabling and re-enabling a core
- Launch an application, first thread created on CPU 1
- First thread will stay on CPU 1, next threads will be created on CPU 1 (default Linux)
- All the threads will be on CPU 1 forever!

In addition to this, when domains re-built, levels 3 and 4 not re-built...

- I.e., no balancing between directly connected or 1-hop CPUs (i.e. any CPU)
- Happens for instance when disabling and re-enabling a core
- Launch an application, first thread created on CPU 1
 - First thread will stay on CPU 1, next threads will be created on CPU 1 (default Linux)
- All the threads will be on CPU 1 forever!



In addition to this, when domains re-built, levels 3 and 4 not re-built...

- I.e., no balancing between directly connected or 1-hop CPUs (i.e. any CPU)
- Happens for instance when disabling and re-enabling a core
- Launch an application, first thread created on CPU 1
 - First thread will stay on CPU 1, next threads will be created on CPU 1 (default Linux)
- All the threads will be on CPU 1 forever!

With bug	After fix	Improvement
122	23	5.2x
134	5.4	25x
72	18	4x
2196	16	137x
	With bug 122 134 72 2196	With bug After fix 122 23 134 5.4 72 18 2196 16





Number of threads in run queue: 0

• Until now, we analyzed the behavior of the the periodic, (buggy) hierarchical load balancing that uses (buggy) scheduling domains

• Until now, we analyzed the behavior of the the periodic, (buggy) hierarchical load balancing that uses (buggy) scheduling domains

But there is another way load is balanced: threads get to pick on which core they get woken up when they are done blocking (after a lock acquisition, an I/O)...

• Until now, we analyzed the behavior of the the periodic, (buggy) hierarchical load balancing that uses (buggy) scheduling domains

- But there is another way load is balanced: threads get to pick on which core they get woken up when they are done blocking (after a lock acquisition, an I/O)...
- Here is how it works: when a thread wakes up, it looks for non-busy cores on the same CPU in order to decide on which core it should wake up.

• Until now, we analyzed the behavior of the the periodic, (buggy) hierarchical load balancing that uses (buggy) scheduling domains

- But there is another way load is balanced: threads get to pick on which core they get woken up when they are done blocking (after a lock acquisition, an I/O)...
- Here is how it works: when a thread wakes up, it looks for non-busy cores on the same CPU in order to decide on which core it should wake up.
- Only cores that are on the same CPU, in order to improve data locality...

• Until now, we analyzed the behavior of the the periodic, (buggy) hierarchical load balancing that uses (buggy) scheduling domains

• But there is another way load is balanced: threads get to pick on which core they get woken up when they are done blocking (after a lock acquisition, an I/O)...

• Here is how it works: when a thread wakes up, it looks for non-busy cores on the same CPU in order to decide on which core it should wake up.

• Only cores that are on the same CPU, in order to improve data locality...

Wait, does that work?

Commercial DB with TPC-H, 64 threads on 64 cores, nothing else on the machine.

- Commercial DB with TPC-H, 64 threads on 64 cores, nothing else on the machine.
- With threads pinned to cores, works fine. With Linux scheduling, execution much slower, phases with overloaded cores while there are long-term idle cores!

- Commercial DB with TPC-H, 64 threads on 64 cores, nothing else on the machine.
- With threads pinned to cores, works fine. With Linux scheduling, execution much slower, phases with overloaded cores while there are long-term idle cores!



Slowed down execution

- Commercial DB with TPC-H, 64 threads on 64 cores, nothing else on the machine.
- With threads pinned to cores, works fine. With Linux scheduling, execution much slower, phases with overloaded cores while there are long-term idle cores!



Slowed down execution





Beginning: 8 threads / CPU, cores busy

 Occasionally, 1 DB thread migrated to other CPU because transient thread appeared during rebalancing which looked like imbalance (only instant loads considered)



- Occasionally, 1 DB thread migrated to other CPU because transient thread appeared during rebalancing which looked like imbalance (only instant loads considered)
- Now, 9 threads on one CPU, and 7 on another one. CPU with 9 threads slow, slows down all execution because all threads wait for each other (barriers), i.e. idle cores everywhere...



- Occasionally, 1 DB thread migrated to other CPU because transient thread appeared during rebalancing which looked like imbalance (only instant loads considered)
- Now, 9 threads on one CPU, and 7 on another one. CPU with 9 threads slow, slows down all execution because all threads wait for each other (barriers), i.e. idle cores everywhere...
- Barriers: threads keep sleeping and waking up, but extra thread never wakes up on idle core, because waking up algorithm only considers local CPU!



- Occasionally, 1 DB thread migrated to other CPU because transient thread appeared during rebalancing which looked like imbalance (only instant loads considered)
- Now, 9 threads on one CPU, and 7 on another one. CPU with 9 threads slow, slows down all execution because all threads wait for each other (barriers), i.e. idle cores everywhere...
- Barriers: threads keep sleeping and waking up, but extra thread never wakes up on idle core, because waking up algorithm only considers local CPU!
- Periodic rebalancing can't rebalance load most of the time because many idle cores
 ⇒ Hard to see an imbalance between 9-thread and 7-thread CPU...



- Occasionally, 1 DB thread migrated to other CPU because transient thread appeared during rebalancing which looked like imbalance (only instant loads considered)
- Now, 9 threads on one CPU, and 7 on another one. CPU with 9 threads slow, slows down all execution because all threads wait for each other (barriers), i.e. idle cores everywhere...
- Barriers: threads keep sleeping and waking up, but extra thread never wakes up on idle core, because waking up algorithm only considers local CPU!
- Periodic rebalancing can't rebalance load most of the time because many idle cores
 ⇒ Hard to see an imbalance between 9-thread and 7-thread CPU...
- "Solution": wake up on core idle for the longest time (not great for energy)



- Periodic rebalancing can't rebalance load most of the time because many idle cores \Rightarrow Hard to see an imbalance between 9-thread and 7-thread CPU...
- "Solution": wake up on core idle for the longest time (not great for energy)

Load balancing on a multicore machine usually considered a solved problem

- Load balancing on a multicore machine usually considered a solved problem
- To recap, on Linux, load balancing works that way:
 - Hierarchical rebalancing uses a metric named load,

- Load balancing on a multicore machine usually considered a solved problem
- To recap, on Linux, load balancing works that way:
 - Hierarchical rebalancing uses a metric named load,
 - to periodically balance threads between scheduling domains.

- Load balancing on a multicore machine usually considered a solved problem
- To recap, on Linux, load balancing works that way:
 - Hierarchical rebalancing uses a metric named load,
 - to periodically balance threads between scheduling domains.
 - In addition to this, threads balance load by selecting core where to wake up.

- Load balancing on a multicore machine usually considered a solved problem
- To recap, on Linux, load balancing works that way:
 - Hierarchical rebalancing uses a metric named load,
 - ↑ Found fundamental issue here
 - to periodically balance threads between scheduling domains.
 - In addition to this, threads balance load by selecting core where to wake up.

- Load balancing on a multicore machine usually considered a solved problem
- To recap, on Linux, load balancing works that way:
 - Hierarchical rebalancing uses a metric named load,
 - **† Found fundamental issue here**
 - to periodically balance threads between scheduling domains.
 - **† Found fundamental issue here**
 - In addition to this, threads balance load by selecting core where to wake up.

- Load balancing on a multicore machine usually considered a solved problem
- To recap, on Linux, load balancing works that way:
 - Hierarchical rebalancing uses a metric named load,
 - **† Found fundamental issue here**
 - to periodically balance threads between scheduling domains.
 - ↑ Found fundamental issue here
 - In addition to this, threads balance load by selecting core where to wake up.
 - ↑ Found fundamental issue here

- Load balancing on a multicore machine usually considered a solved problem
- To recap, on Linux, load balancing works that way:
 - Hierarchical rebalancing uses a metric named load,
 - **† Found fundamental issue here**
 - to periodically balance threads between scheduling domains.
 - **†** Found fundamental issue here
 - In addition to this, threads balance load by selecting core where to wake up.
 - ↑ Found fundamental issue here

Wait, was anything working at all? 🙂

Many major issues went unnoticed for years in the scheduler... How can we prevent this from happening again?

Many major issues went unnoticed for years in the scheduler... How can we prevent this from happening again?

Code testing

- No clear fault (no crash, no deadlock, etc.)
- Existing tools don't target these bugs

Many major issues went unnoticed for years in the scheduler... How can we prevent this from happening again?

Code testing

- No clear fault (no crash, no deadlock, etc.)
- Existing tools don't target these bugs

Performance regression

- Usually done with 1 app on a machine to avoid interactions
- Insufficient coverage

Many major issues went unnoticed for years in the scheduler... How can we prevent this from happening again?

Code testing

- No clear fault (no crash, no deadlock, etc.)
- Existing tools don't target these bugs

Performance regression

- Usually done with 1 app on a machine to avoid interactions
- Insufficient coverage

Model checking, formal proofs

Complex, parallel code: so far, nobody knows how to do it...

Idea 1: short-term hack — implemented a sanity checker

Idea 1: short-term hack — implemented a sanity checker

Idea 1: short-term hack — implemented a sanity checker



Idea 1: short-term hack — implemented a sanity checker



Idea 2: fine-grained tracers!

Idea 2: fine-grained tracers!

- Built a simple one, turned out to be the only way to really understand what happens
- Aggregate metrics (CPI, cache misses, etc.) not precise enough



Idea 2: fine-grained tracers!

- Built a simple one, turned out to be the only way to really understand what happens
- Aggregate metrics (CPI, cache misses, etc.) not precise enough



Could really be improved!

Idea 3: produce a dedicated profiler!

- Idea 3: produce a dedicated profiler!
 - Lack of tools!

- Idea 3: produce a dedicated profiler!
 - Lack of tools!
 - Possible to detect if slowdown comes from scheduler or application?
 - Would avoid a lot of wasted time!

- Idea 3: produce a dedicated profiler!
 - Lack of tools!
 - Possible to detect if slowdown comes from scheduler or application?
 - Would avoid a lot of wasted time!
 - Follow threads, and see if often on overloaded cores when shouldn't have?
 - Detect if threads unnecessarily moved to core/node that leads to many cache misses?

Idea 4: produce good scheduler benchmarks!

- Idea 4: produce good scheduler benchmarks!
 - Really needed, and virtually inexistent!

- Idea 4: produce good scheduler benchmarks!
 - Really needed, and virtually inexistent!
 - Not an easy problem: insane coverage needed!
 - Using combination of many real applications: configuration nightmare!

- Idea 4: produce good scheduler benchmarks!
 - Really needed, and virtually inexistent!
 - Not an easy problem: insane coverage needed!
 - Using combination of many real applications: configuration nightmare!
 - Simulated workloads?
 - Have to do elaborate work: spinning and sleeping not efficient
 - Have to be representative of reality, have to cover corner cases
 - Use machine learning? Genetic algorithms?

Idea 5: switch to simpler schedulers, easier to reason about!

Let's take a step back: *why* did we end up in this situation?

- Idea 5: switch to simpler schedulers, easier to reason about!
- Let's take a step back: *why* did we end up in this situation?
- Linux used for many classes of applications (big data, n-tier, cloud, interactive, DB, HPC...)

Idea 5: switch to simpler schedulers, easier to reason about!

Let's take a step back: *why* did we end up in this situation?

- Linux used for many classes of applications (big data, n-tier, cloud, interactive, DB, HPC...)
- Multicore architectures increasingly diverse and complex!

Idea 5: switch to simpler schedulers, easier to reason about!

Let's take a step back: *why* did we end up in this situation?

- Linux used for many classes of applications (big data, n-tier, cloud, interactive, DB, HPC...)
- Multicore architectures increasingly diverse and complex!
- **Result:** very complex monolithic scheduler supposed to work in all situations!
 - Many heuristics interact in complex, unpredictable ways
 - Some features greatly complexify, e.g., load balancing (tasksets, cgroups/autogroups...)

Idea 5: switch to simpler schedulers, easier to reason about!

Let's take a step back: *why* did we end up in this situation?

- Linux used for many classes of applications (big data, n-tier, cloud, interactive, DB, HPC...)
- Multicore architectures increasingly diverse and complex!
- **Result:** very complex monolithic scheduler supposed to work in all situations!
 - Many heuristics interact in complex, unpredictable ways
 - Some features greatly complexify, e.g., load balancing (tasksets, cgroups/autogroups...)
- Keeps getting worse!
 - **E.g., task_struct:** 163 fields in Linux 3.0 (07/2011), 215 fields in 4.6 (05/2016)
 - **20,000 lines of C!**

Idea 5: switch to simpler schedulers, easier to reason about!



Idea 5: switch to simpler schedulers, easier to reason about!

Proving the scheduler implementation correct: not doable!

- Way too much code for current technology
- We'd need to detect high-level abstractions from low-level C: a challenge!
- Even if we managed that, how do we keep up with updates?
 - Code keeps evolving with new architectures and application needs...

Idea 5: switch to simpler schedulers, easier to reason about!

Proving the scheduler implementation correct: not doable!

- Way too much code for current technology
- We'd need to detect high-level abstractions from low-level C: a challenge!
- Even if we managed that, how do we keep up with updates?
 - Code keeps evolving with new architectures and application needs...
- We need another approach...

- Idea 5: switch to simpler schedulers, easier to reason about!
- Write simple, schedulers with proven properties !
 - A scheduler is tailored to a (class of) parallel application(s)
 - Specific thread election criterion, load balancing criterion, state machine with events...
 - Machine partitioned into sets of cores that run ≠ schedulers
 - Scheduler deployed together with (an) application(s) on a partition

- Idea 5: switch to simpler schedulers, easier to reason about!
- Write simple, schedulers with proven properties !
 - A scheduler is tailored to a (class of) parallel application(s)
 - Specific thread election criterion, load balancing criterion, state machine with events...
 - Machine partitioned into sets of cores that run ≠ schedulers
 - Scheduler deployed together with (an) application(s) on a partition
- Through a DSL, for two reasons:
 - Much easier, safer and less bug-prone than writing low-level C kernel code !
 - Clear abstractions, possible to reason about them and prove properties
 - Work conservation, load balancing live and in finite # or rounds, valid hierarchy...

Towards Proving Optimistic Multicore Schedulers

Baptiste Lepers, Willy Zwaenepoel EPFL first.last@epfl.ch

Jean-Pierre Lozi Université Nice Sophia-Antipolis

jplozi@unice.fr

Nicolas Palix Université Grenoble Alpes nicolas.palix@univ-grenoblealpes.fr

Redha Gouicem, Julien Sopena, Julia Lawall, Gilles Muller

Sorbonne Universités, Inria, LIP6 first.last@lip6.fr



WHERE DO WE GO FROM HERE?

- Idea 5: switch to simpler schedulers, easier to reason about!
- Write simple, schedulers with proven properties !
 - A scheduler is tailored to a (class of) parallel application(s)
 - Specific thread election criterion, load balancing criterion, state machine with events...
 - Machine partitioned into sets of cores that run ≠ schedulers
 - Scheduler deployed together with (an) application(s) on a partition
- Through a DSL, for two reasons:
 - Much easier, safer and less bug-prone than writing low-level C kernel code !
 - Clear abstractions, possible to reason about them and prove properties
 - Work conservation, load balancing live and in finite # or rounds, valid hierarchy...

- Idea 6: 333
- Any other ideas?

Scheduling (as in dividing CPU cycles among theads) was thought to be a solved problem.

- Scheduling (as in dividing CPU cycles among theads) was thought to be a solved problem.
- Analysis: fundamental issues in the load metric, scheduling domains, scheduling choices...

- Scheduling (as in dividing CPU cycles among theads) was thought to be a solved problem.
- Analysis: fundamental issues in the load metric, scheduling domains, scheduling choices...
- Very bug-prone implementation following years of adapting to hardware

- Scheduling (as in dividing CPU cycles among theads) was thought to be a solved problem.
- Analysis: fundamental issues in the load metric, scheduling domains, scheduling choices...
- Very bug-prone implementation following years of adapting to hardware
- Can't ensure simple "invariant": no idle cores while overloaded cores

- Scheduling (as in dividing CPU cycles among theads) was thought to be a solved problem.
- Analysis: fundamental issues in the load metric, scheduling domains, scheduling choices...
- Very bug-prone implementation following years of adapting to hardware
- Can't ensure simple "invariant": no idle cores while overloaded cores
- Proposed fixes: not always satisfactory

- Scheduling (as in dividing CPU cycles among theads) was thought to be a solved problem.
- Analysis: fundamental issues in the load metric, scheduling domains, scheduling choices...
- Very bug-prone implementation following years of adapting to hardware
- Can't ensure simple "invariant": no idle cores while overloaded cores
- Proposed fixes: not always satisfactory
- What can we do? Many things to explore!

- Scheduling (as in dividing CPU cycles among theads) was thought to be a solved problem.
- Analysis: fundamental issues in the load metric, scheduling domains, scheduling choices...
- Very bug-prone implementation following years of adapting to hardware
- Can't ensure simple "invariant": no idle cores while overloaded cores
- Proposed fixes: not always satisfactory
- What can we do? Many things to explore!

Our takeaway: more research must be directed towards implementing efficient and reliable schedulers for multicore architectures!

- Scheduling (as in dividing CPU cycles among theads) was thought to be a solved problem.
- Analysis: fundamental issues in the load metric, scheduling domains, scheduling choices...
- Very bug-prone implementation following years of adapting to hardware
- Can't ensure simple "invariant": no idle cores while overloaded cores
- Proposed fixes: not always satisfactory
- What can we do? Many things to explore!

Our takeaway: more research must be directed towards implementing efficient and reliable schedulers for multicore architectures!
Your turn!