	License
	Copyright © 2008 Ciaran McHale.
Multi-threaded Performance Pitfalls Ciaran McHale	Permission is hereby granted, free of charge, to any person obtaining a copy of this training course and associated documentation files (the "Training Course"), to deal in the Training Course without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Training Course, and to permit persons to whom the Training Course is furnished to do so, subject to the following conditions:
	The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Training Course.
	THE TRAINING COURSE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE TRAINING COURSE OR THE USE OR OTHER DEALINGS IN THE TRAINING COURSE.
CiaranMcHale.com Complexity explained simply	1 Multi-threaded Performance Pitfalls 2
Purpose of this presentation	
Some issues in multi-threading are counter-intuitive	
 Ignorance of these issues can result in poor performance Performance can actually get <i>worse</i> when you add more CPUs 	
This presentation explains the counter-intuitive issues	1. A case study
	3

Architectural diagram	Architectural notes
web browser bouter bload balancing router J2EE App Server1 CORBA C++ server on 8-CPU Solaris box J2EE App Server2 J2EE Server2	 The customer felt J2EE was slower than CORBA/C++ So, the architecture had: Multiple J2EE App Servers acting as clients to Just one CORBA/C++ server that ran on an 8-CPU Solaris box The customer assumed the CORBA/C++ server "should be able to cope with the load"
Multi-threaded Performance Pitfalls 5 Strange problems were observed	Multi-threaded Performance Pitfalls 6
 Throughput of the CORBA server <i>decreased</i> as the number of CPUs increased It ran fastest on 1 CPU It ran slower but "fast enough" with moderate load on 4 CPUs (development machines) It ran very slowly on 8 CPUs (production machine) The CORBA server ran faster if a thread pool limit was imposed 	2. Analysis of the problems
 Under a high load in production: Most requests were processed in < 0.3 second But some took up to a minute to be processed A few took up to 30 minutes to be processed This is <i>not</i> what you hope to see Multi-threaded Performance Pitfalls 	ε

Multi-threaded Performance Pitfalls

What went wrong?	Cache consistency
Investigation showed that scalability problems were caused by	RAM access is much slower than speed of CPU
a combination of:	- Solution: high-speed cache memory sits between CPU and RAM
- Cache consistency in multi-CPU machines	Cache memory works great:
- Unfair mutex wakeup semantics	- In a single-CPU machine
These issues are discussed in the following slides	 In a multi-CPU machine if the threads of a process are "bound" to a CPU
Another issue contributed (slightly) to scalability problems:	 Cache memory can backfire if the threads in a program are spread over all the CPUs:
- Bottlenecks in application code	- Each CPU has a separate cache
- A discussion of this is outside the scope of this presentation	 Cache consistency protocol require cache flushes to RAM (cache consistency protocol is driven by calls to lock() and unlock())
Multi-threaded Performance Pitfalls 9 Cache consistency (cont')	Multi-threaded Performance Pitfalls Unfair mutex wakeup semantics
- Overhead of eache consistency protocole werease as	- A residence net assertate First In First Out (FIFO) suchas
Overhead of cache consistency protocols worsens as:	 A mutex does not guarantee First In First Out (FIFO) wakeup semantics
- Overhead of a cache synchronization increases	semantics
	• • • •
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) Frequency of cache synchronization increases 	 semantics To do so would prevent two important optimizations (discussed on the following slides)
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides:
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) Frequency of cache synchronization increases (this increases with the rate of mutex lock() and unlock() calls) 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides: Unfair wakeup semantics
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) Frequency of cache synchronization increases (this increases with the rate of mutex lock() and unlock() calls) Lessons: 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides: Unfair wakeup semantics Can cause <i>temporary</i> starvation of a thread
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) Frequency of cache synchronization increases (this increases with the rate of mutex lock() and unlock() calls) 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides: Unfair wakeup semantics Can cause <i>temporary</i> starvation of a thread But guarantees to avoid <i>infinite</i> starvation
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) Frequency of cache synchronization increases (this increases with the rate of mutex lock() and unlock() calls) Lessons: Increasing number of CPUs can <i>decrease</i> performance of a server 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides: Unfair wakeup semantics Can cause <i>temporary</i> starvation of a thread
 (this increases as the number of CPUs increase) Frequency of cache synchronization increases (this increases with the rate of mutex lock() and unlock() calls) Lessons: Increasing number of CPUs can <i>decrease</i> performance of a server Work around this by: 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides: Unfair wakeup semantics Can cause <i>temporary</i> starvation of a thread But guarantees to avoid <i>infinite</i> starvation
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) Frequency of cache synchronization increases (this increases with the rate of mutex lock() and unlock() calls) Lessons: Increasing number of CPUs can <i>decrease</i> performance of a server Work around this by: Having multiple server processes instead of just one Binding each process to a CPU (avoids need for cache 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides: Unfair wakeup semantics Can cause <i>temporary</i> starvation of a thread But guarantees to avoid <i>infinite</i> starvation
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) Frequency of cache synchronization increases (this increases with the rate of mutex lock() and unlock() calls) Lessons: Increasing number of CPUs can <i>decrease</i> performance of a server Work around this by: Having multiple server processes instead of just one Binding each process to a CPU (avoids need for cache synchronization) 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides: Unfair wakeup semantics Can cause <i>temporary</i> starvation of a thread But guarantees to avoid <i>infinite</i> starvation
 Overhead of a cache synchronization increases (this increases as the number of CPUs increase) Frequency of cache synchronization increases (this increases with the rate of mutex lock() and unlock() calls) Lessons: Increasing number of CPUs can <i>decrease</i> performance of a server Work around this by: Having multiple server processes instead of just one Binding each process to a CPU (avoids need for cache synchronization) Try to minimize need for mutex lock() and unlock() in application 	 semantics To do so would prevent two important optimizations (discussed on the following slides) Instead, a mutex provides: Unfair wakeup semantics Can cause <i>temporary</i> starvation of a thread But guarantees to avoid <i>infinite</i> starvation

Multi-threaded Performance Pitfalls

Unfair mutex wakeup semantics (cont')	Mutex optimization 1
 Why does a mutex <i>not</i> provide fair wakeup semantics? Because most of the time, speed matter more than fairness When FIFO wakeup semantics are required, developers can write a FIFOMutex class and take a performance hit 	 Pseudo-code: void lock()
Multi-threaded Performance Pitfalls 13 Mutex optimization 2	Multi-threaded Performance Pitfalls
 Assume several threads concurrently execute the following code: for (i = 0; i < 1000; i++) { lock(a_mutex); process(data[i]); unlock(a_mutex); } A thread context switch is (relatively) expensive Context switching on every unlock() would add a lot of overhead 	3. Improving Throughput
 Solution (this is an unfair optimization): Defer context switches until the end of the current thread's time slice Current thread can <i>repeatedly</i> lock() and unlock() mutex in a single time slice Multi-threaded Performance Pitfalls 15 	1

Improving throughput	
20X increase in throughput was obtained by combination of:	
 Limiting size of the CORBA server's thread pool This Decreased the maximum length of the mutex wakeup queue Which decreased the maximum wakeup time Using several server processes (each with a small thread pool) rather than one server process (with a very large thread pool) Binding each server process to one CPU This avoided the overhead of cache consistency Binding was achieved with the pbind command on Solaris Windows has an equivalent of process binding: 	4. Finishing up
- Use the SetProcessAffinityMask() system call - Or, in Task Manager, right click on a process and choose the menu option (this menu option is visible only if you have a multi-CPU machine) Multi-threaded Performance Pitfalls 17 Recap: architectural diagram	The case study is not an isolated incident
web browser bload balancing router	 The project's high-level architecture is quite common: Multi-threaded clients communicate with a multi-threaded server Server process is not "bound" to a single CPU Server's thread pool size is unlimited (this is the default case in many middleware products) Likely that many projects have similar scalability problems: But the system load is not high enough (yet) to trigger problems Problems are not specific to CORBA They are independent of your choice of middleware technology Multi-core CPUs are becoming more common So, expect to see these scalability issues occurring more frequently
Multi-threaded Performance Pitfalls 19	Multi-threaded Performance Pitfalls 2

Multi-threaded Performance Pitfalls

Summary: important things to remember

- Recognize danger signs:
 - Performance drops as number of CPUs increases
 - Wide variation in response times with a high number of threads
- Good advice for multi-threaded servers:
 - Put a limit on the size of a server's thread pool
 - Have several server processes with a small number of threads instead of one process with many threads
 - Bind each a server process to a CPU

Acknowledgements:

- Ciaran McHale's employer, IONA Technologies (www.iona.com) generously gave permission for this presentation to be released under the stated open-source license.

Multi-threaded Performance Pitfalls

21