
Java Programming in a Multicore World

Angelika Langer
Trainer/Consultant

<http://www.AngelikaLanger.com>

objective

- take look at current trends in concurrent programming
- explain the Java Memory Model
- discuss future trends such as lock-free programming and transactional memory

speaker's relationship to topic

- independent trainer / consultant / author
 - teaching C++ and Java for 10+ years
 - curriculum of a dozen challenging courses
 - co-author of "Effective Java" column
 - author of Java Generics FAQ online
 - Java champion since 2005



© Copyright 2003-2008 by Angelika Langer & Klaus Krefl. All Rights Reserved.
<http://www.AngelikaLanger.com>
last update: 6/14/2008, 11:30

3

agenda

- **concurrency trends**
- synchronization and memory model
- fight the serialization – improve scalability
- future trends

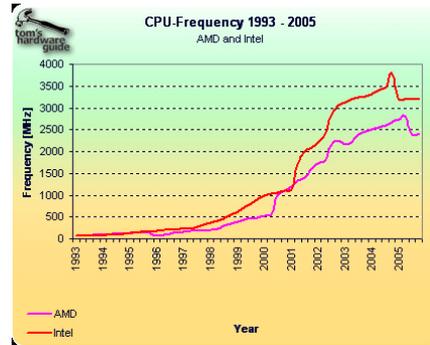


© Copyright 2003-2008 by Angelika Langer & Klaus Krefl. All Rights Reserved.
<http://www.AngelikaLanger.com>
last update: 6/14/2008, 11:30

4

CPU development

- Moore's law:
 - number of transistors doubles every two years
 - since 2004: more cores
 - until 2004: faster ones
 - main reason: heat
- 2 cores became standard 2007
 - 6-12 in 2009 (AMD)
- more complex caches
 - hierarchy



CPU development implies

- new CPU will not solve your performance problems
 - if your program does not scale (well) to multiple cores
 - i.e.: find (and fight) the serialization
- existing programs
 - undetected errors might pop up
 - multi-core + caching uncovers synchronization problems
- Java environment
 - more and more complex work for
 - the byte code compiler, and
 - the JIT compiler

Java history – initial MT support

- mostly built into the language (not into the library)
 - `synchronized` block/method – lock in every object
 - `Object.wait()`, `Object.notify()` – condition in every object
 - ...
- mainly low level functionality
 - no thread pool, no blocking queue, ...
- memory model
 - chapter 17 of the Java Language Specification: Threads and Locks
 - hard to understand,
 - incomplete,
 - violated by JVM implementations



© Copyright 2003-2008 by Angelika Langer & Klaus Krefl. All Rights Reserved.
<http://www.AngelikaLanger.com>
last update: 6/14/2008, 11:30

7

Java history – JDK 5.0 MT support

- rework of existing locks and conditions
 - into the library: `java.util.concurrent.locks`
 - extended functionality
 - timeout for existing locks
 - new locks: read-write-lock
 - approach changed: library is more flexible than language
 - think of C
- high-level abstractions
 - thread pool: `ThreadPoolExecutor`, ...
 - synchronizers: `BlockingQueue`, `CyclicBarrier`, ...
 - support for asynchronous programming: `Future`, ...



© Copyright 2003-2008 by Angelika Langer & Klaus Krefl. All Rights Reserved.
<http://www.AngelikaLanger.com>
last update: 6/14/2008, 11:30

8

(cont.)

- support for lock free programming
 - low-level
 - abstractions from `java.util.concurrent.atomic`
 - high-level
 - 'concurrent' collections: `ConcurrentHashMap`, ...
- reworked memory model
 - cleared up what `volatile` and `final` mean in a MT context
 - defines requirements regarding atomicity, visibility and ordering of operations

Java keeps up to ...

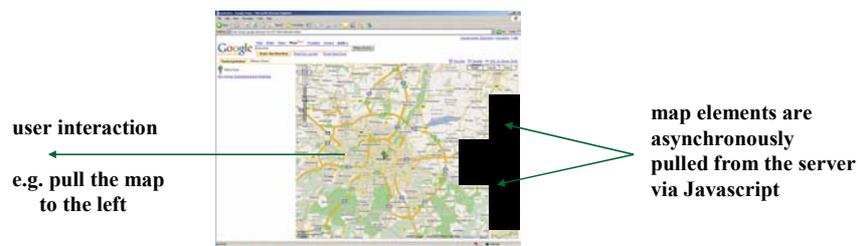
- ... the needs and requirements of the changing MT uses
- more people build MT programs
 - MT patterns and idioms become common knowledge
 - need for high-level abstractions
- more people build Java MT programs
for multiprocessor platforms
 - need for clear and exact memory model
 - wish for better scaling MT abstractions
 - need for lock free programming
- former niche becomes main stream with multi core CPUs

system architecture - current trends

- trend to asynchronous and parallel computing
to **increase throughput**
- Java examples
 - asynchronous I/O
 - 1.4 socket, 5.0 sockets + SSL, 7.0 sockets + file system
 - essential: frees you from one thread per socket
 - but: program structure gets more complex and technical
 - JMS introduced 2001
 - much later than RMI which was part of Java from the beginning
 - effect: EJB became *message-driven beans*

current trends

- example: AJAX (**A**synchronous **J**avascript and **X**ML)
 - means: **decouple** user interaction from HTTP requests
 - traditionally:
 - you select a link / push a button / etc. , and
 - a new page gets loaded into your browser
 - AJAX example: Google Maps



more AJAX

- more asynchronicity: **HTTP push** via Ajax
 - signal an asynchronous event in the browser
 - e.g. incoming telephone call
- solutions boil down to a 'long-lived' HTTP request
 - persistent communication / **long polling** / hybrid polling:
 - request lives, until the event occurs or 'long' timeout occurs (5-10min)
 - event (or timeout) is signaled in the response
 - new request to poll the next event
 - comet style / **HTTP streaming**:
 - request lives, until the client goes away
 - all data is send from the server to the client in the same response

a small problem ...

- traditional servlet programming
 - one thread:
 - receives HTTP request
 - determines what has to be done
 - gathers the data (and renders the new page)
 - sends all this back to the client in a response
- what about an long-lived open HTTP request ?
 - that waits for an external event,
 - e.g. the incoming telephone call
- allocates a thread until
the event occurs / client goes away !
 - with 50000 users on the server ?!?!?

asynchronous web servers

- decouple the request from the response
- Jetty 6 Continuation
 - `Continuation.suspend()`, `Continuation.resume()`
 - <http://docs.codehaus.org/display/JETTY/Continuations>
- Tomcat 6.0 Comet
 - Comet module allows to process I/O asynchronously
 - <http://tomcat.apache.org/tomcat-6.0-doc/aio.html>
- Java standard for asynchronous web server
 - JSR 315 = Servlet 3.0 specification
 - scheduled finish by the end of 2008
- underlying concept: asynchronous I/O

more and more asynchronicity

- not only web server – other servers too
- SOA (service oriented architecture)
 - service -> service -> service ...
 - you don't want to have a waiting thread in each of the server
 - i.e.
 - asynchronous handling of the request
 - MOM (message oriented middleware), means often JMS in Java
- bottom line:
 - you need **multiple threads** and some **synchronization** of these to tie the external asynchronous channels to your program

agenda

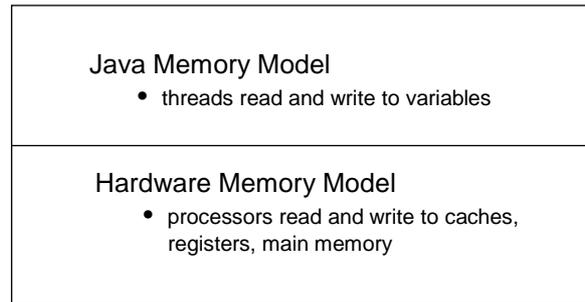
- concurrency trends
- **synchronization and memory model**
- fight the serialization – improve scalability
- future trends

motivation - why does JMM matter?

- JMM = **J**ava **M**emory **M**odel
- understanding JMM reveals errors in existing programs
 - undetected errors might pop up
 - multi-core + caching uncovers synchronization problems

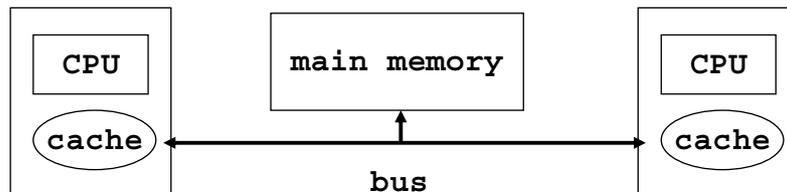
Java Memory Model (JMM)

- specifies minimal guarantees given by the JVM
 - about when writes to variables become visible to other threads
- is an abstraction on top of hardware memory models



Java memory model

- JMM resembles abstract SMP (symmetric multi processing) machine
- key ideas:
 - all threads share the main memory
 - each thread uses a local working memory
 - flushing or refreshing working memory to/from main memory must comply to JMM rules



Java memory model

JMM rules address 3 intertwined issues:

- atomicity
 - which operations must have indivisible effects ?
- visibility
 - under which conditions are the effects of operations taken by one thread visible to other threads ?
- ordering
 - under which conditions can the effects of operations appear out of order to any given thread ?

"operations" means:

- reads and writes to memory cells representing Java variables

JMM in practice

- examples:

atomicity

- access to variables of primitive type (except long/double) are atomic
- execution of operations in a synchronized block is atomic

visibility

- values written to a volatile variable are visible to other threads

ordering

- effects of operations in a synchronized block appear in order
- accesses to volatile variables appear in order

sequential consistency

- **sequential consistency** is a convenient (yet unrealistic) mental model:
 - imagine a single order of execution of all programm operations (regardless of the processors used)
 - each read of a variable will see the last write in the execution order
- JMM does NOT guarantee sequential consistency
 - reordering is generally permitted
 - specific rules for synchronization, thread begin/end, volatile and final variables

agenda

- concurrency trends
- synchronization and memory model
 - atomicity
 - **visibility**
 - ordering
- fight the serialization – improve scalability
- future trends

need for visibility

```
private int[] array;
private int cnt = 0; ← must be volatile
...
public synchronized void push(int elm) { array[cnt++] = elm; }
public synchronized int pop()         { return(array[--cnt]); }
public int size()                       { return cnt; }
...
```

- access to cnt is atomic
 - no synchronization in size() needed
- visibility problem
 - writes performed in one thread need not be visible to other threads
 - i.e. modification of cnt in push()/pop() need not be visible to size()
- volatile is needed not for atomicity, but for visibility

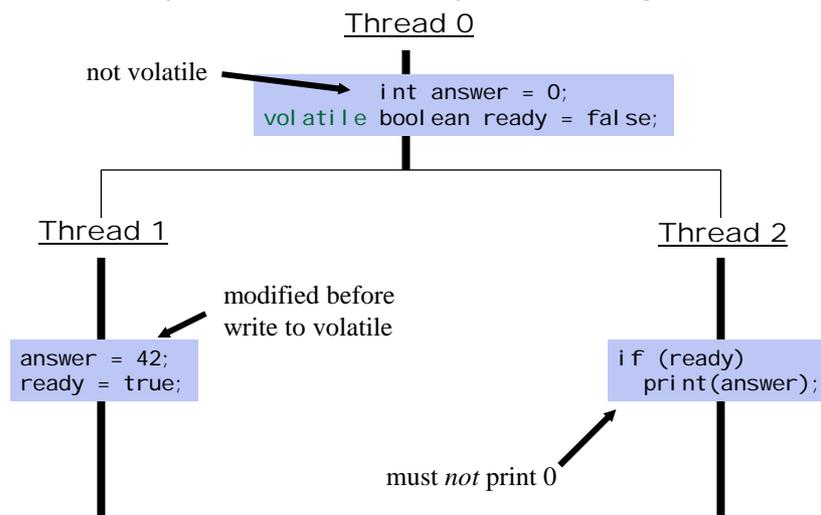
visibility guarantees

- changes made in one thread are guaranteed to be visible to other threads under the following conditions:
 - explicit synchronization
 - thread start and termination
 - read / write of volatiles
 - first read of finals

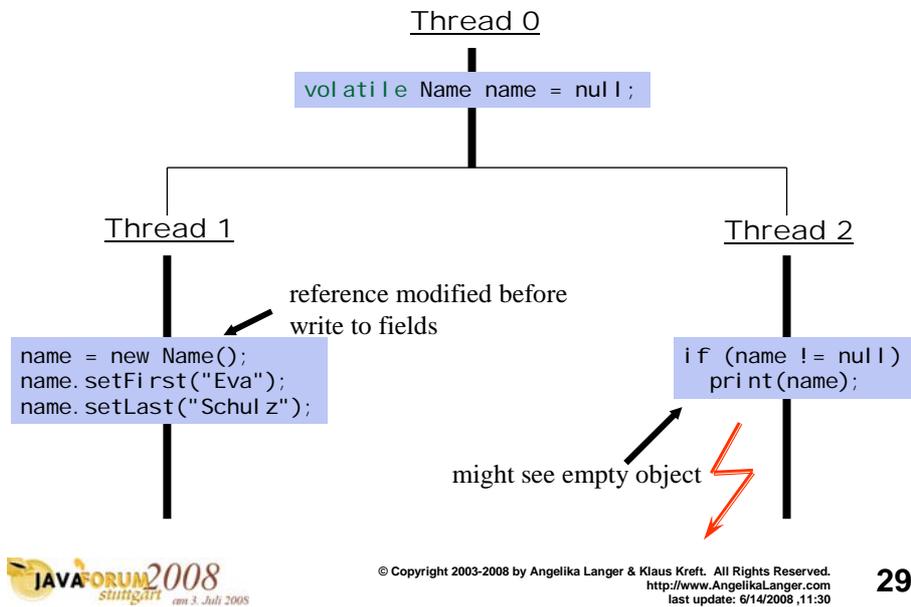
visibility guarantee: read / write of volatiles

- reading a volatile forces a reload from main memory
- writing to a volatile forces a flush to main memory
- matches our expectation
 - when a thread reads a volatile, then all writes are visible that any other thread performed prior to a write to the same volatile
- how about volatile references ... ?
- volatile is *not* transitive
 - read/write of a volatile reference affects the reference, but not the referenced object (or array)

volatile (since Java 5) - example



non-transitive volatile - example



volatile references

- what do we do to also make the modified object visible?
 - make all fields of referenced object volatile
 - problem for arrays: array elements cannot be declared volatile
 - modify elements before assignment to volatile reference
 - all changes made prior to writing to the volatile are flushed
 - use explicit synchronization
 - viable fallback, at the expense of synchronization overhead

final vs. volatile vs. synchronization

- explicit synchronization is expensive
- volatile is cheaper but still relatively expensive
 - due to the need of memory barriers
- final is even cheaper
 - memory barrier is needed only once (after construction)

agenda

- concurrency trends
- synchronization and memory model
- **fight the serialization – improve scalability**
- future trends

Amdahl's law

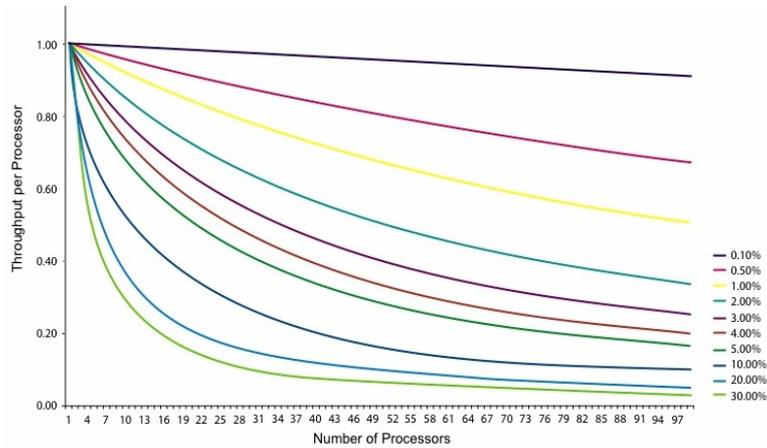
- named after computer architect Gene Amdahl
 - "Validity of the Single Processor Approach to Achieving Large-Scale Computing Capabilities", AFIPS Conference Proceedings, (30), pp. 483-485, 1967.
 - Gene Amdahl has approved the use of his complete text in the Usenet comp.sys.super news group FAQ which goes out on the 20th of each month
- used in parallel computing to predict the theoretical maximum speedup using multiple processors

(cont.)

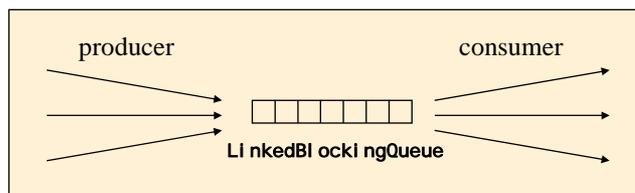
- idea: divide work into **serial** and **parallel** portions
 - serial work cannot be sped up by adding resources
 - parallelizable work can
- Amdahl's Law: $\text{speedup} \leq \frac{1}{\left(F + \frac{(1 - F)}{N}\right)}$
 - F is the fraction that must be serialized
 - N is the number of CPUs
- with $N \rightarrow \infty$, speedup $\rightarrow 1/F$
 - with 50% serialization,
 - your program can only speed up by a factor of 2 (with: ∞ CPUs)
- naïve idea: from 1 to 2 CPUs = factor of 2 ?

(cont.)

- fight serialization to improve performance



example



- looks highly parallelizable
 - (if producers are slow increase their thread pool)
- 0% serialized ?
 - no!
 - need synchronization to maintain the queue's integrity

LinkedListQueue.offer()

```
public boolean offer(E o) {
    if (o == null) throw new NullPointerException();
    final AtomicInteger count = this.count;
    if (count.get() == capacity)
        return false;
    int c = -1;
    final ReentrantLock putLock = this.putLock;
    putLock.lock();
    try {
        if (count.get() < capacity) {
            insert(o);
            c = count.getAndIncrement();
            if (c + 1 < capacity)
                notFull.signal();
        }
    } finally {
        putLock.unlock();
    }
    if (c == 0)
        signalNotEmpty();
    return c >= 0;
}
```

(cont.)

- Doug Lea did an excellent job with the implementation
 - highly optimized
 - split lock: put / take
 - count guarded lock-free
 - stack-local variables to speed up the execution inside the critical region
 - ...
- structural problem
 - serialization of offering threads (producers)
 - similar serialization of getting threads (consumers)

serialization

- where/when threads demand concurrent access
- often hidden
 - in frameworks / third party abstractions
- other area: asynchronous service architecture
 - example: `java.nio.channels.Selector`
 - section on concurrency in the respective JavaDoc
 - management to send back the result asynchronously
 - Jetty continuation

fight the serialization ...

... try to reduce lock induced serialization

- smallest critical region possible
 - synchronized block vs. synchronized method
 - or use explicit locks
 - speed up execution inside the critical region
 - replace synchronized counters with `AtomicInteger`
- lock splitting / striping
 - guard different state with different locks
 - reduces likelihood of lock contention

fight the serialization ...

... try to eliminate locking entirely

- replace mutable objects with immutable ones
- replace shared objects with thread-local ones
 - e.g. make a copy before passing it to a concurrent thread
- lock-free programming

agenda

- history of concurrency & concurrency trends
- synchronization and memory model
- fight the serialization – improve scalability
- **future trends**

trends

- lock-free programming
 - supported in Java since JDK 5.0
 - `java.util.concurrent.atomic`, and
 - Concurrent collections in `java.util.concurrent`
- transactional memory
 - neither supported in Java nor in any popular programming language at the moment
- commonality
 - avoid locking to avoid serialization

agenda

- history of concurrency & concurrency trends
- synchronization and memory model
- fight the serialization – improve scalability
- **future trends**
 - **lock free programming**
 - transactional memory

CAS

- modern processors have a primitive called CAS (**compare-and-swap**)
- a CAS operation includes three operands
 - a memory location
 - the expected old value
 - a new value
- processor will atomically update the location to the new value
 - if the value that is there matches the expected old value
 - otherwise it will do nothing
 - it returns the value that was at that location prior to the CAS instruction

CAS permits atomic read-modify-write

- CAS allows an algorithm to execute a read-modify-write sequence
 - without fear of another thread modifying variable in meantime
 - if another thread did modify variable, CAS would detect it (and fail)
 - and algorithm could retry operation
- CAS-like operation are available in JDK 5.0 as "atomic variables"
 - based on the underlying system/hardware/CPU support
 - `java.util.concurrent.atomic`

example - thread-safe counter

```
public class SafeCounter {
    private volatile int value;
    public int getValue() { return value; }
    public synchronized int increment() { return ++value; }
    public synchronized int decrement() { return --value; }
}
```

- increment() / decrement() are read-modify-write operations and must be atomic
 - atomic read-modify-write cannot be achieved by making instance variable volatile
 - need to be synchronized
- get() without synchronization, since value is volatile

example - thread-safe counter – lock free

```
public class AtomicCounter {
    private AtomicInteger value;
    public int getValue() { return value.get(); }
    public int increment() {
        int oldValue = value.get();
        while (!value.compareAndSet(oldValue,
                                   oldValue + 1))
            oldValue = value.get();
        return oldValue + 1;
    }
    public int decrement() {
        int oldValue = value.get();
        while (!value.compareAndSet(oldValue,
                                   oldValue - 1))
            oldValue = value.get();
        return oldValue - 1;
    }
}
```

is atomic

lock-free

- advantages
 - fast (~ 4 times faster than best locks)
 - deadlock immunity
 - ...
- disadvantages
 - hard to program !!!
 - no simple straight forward approach as with locks
 - ...

hard to program, but what you can do

- some strategies
 - e.g. lock-free counter, ABA problem, ...
 - no single best resource of information known
 - best to search the web for 'lock free programming'
- algorithms for standard data structures
 - map, linked list, ...
 - Concurrent collections from `java.util.concurrent`
 - use these in your program
 - or these in combination with locks

agenda

- history of concurrency & concurrency trends
- synchronization and memory model
- fight the serialization – improve scalability
- future trends
 - lock free programming
 - **transactional memory**

transactional memory ...

- ... or **software transactional memory** (STM)
- similar to optimistic strategies in database transactions
 - e.g. *optimistic locking pattern for EJBs*
- very intuitive

```
public void addName(String name) {  
    atomic {  
        nameCount++;  
        nameList.add(name);  
    }  
}
```

optimistic strategy

- thread does modifications to shared memory/object
 - without regards what the other threads are doing
- finished modifications
 - commit
 - verification that no other thread made concurrent modifications
 - abort and rollback
 - concurrent modifications occurred
 - error handling: (in most cases) retry of the transactions
- increased concurrency vs. overhead of retrying transactions that failed

evaluation

PRO

- very intuitive, e.g. update an object in shared memory
 - close to the original Java approach of synchronized blocks

CON

- not every operation can be rolled back
 - what about those that are not memory based ?, e.g. unbuffered I/O
- no popular programming language supports STM
 - must use a more experimental language, e.g. Clojure
 - dynamic language, Lisp dialect, compiles to JVM bytecode
 - <http://clojure.sourceforge.net/>

wrap-up

- a trend towards concurrent, asynchronous computing
 - MT initially for better structure
 - today to overcome synchronicity (messaging, AJAX, ...)
- multicore architecture might reveal yet undetected bugs
 - due to memory model issues (atomicity, visibility, ordering)
- multicore architectures need scalable software to be useful
 - avoid serialization - increase concurrency - Amdahl's law
- a gaze into the crystal ball
 - lock-free programming is already in use (by experts)
 - transactional memory might ease concurrent programming some time in the future

authors

Angelika Langer

Training & Mentoring

Object-Oriented Software Development in C++ & Java

Email: contact@AngelikaLanger.com

http: www.AngelikaLanger.com

Klaus Krefl

Siemens Enterprise Communications GmbH & Co. KG, Munich, Germany

Email: klaus.krefl@siemens.com

Java Programming in a Multicore World

Q & A



© Copyright 2003-2008 by Angelika Langer & Klaus Kref. All Rights Reserved.
<http://www.AngelikaLanger.com>
last update: 6/14/2008, 11:30

57