# **Programming Languages**

Concurrency: Memory Consistency

Dr. Michael Petter
Winter term 2019

Thread A

```
void foo(void) {
  a = 1;
  b = 1;
}
```

Thread B

```
void bar(void) {
  while (b == 0){};
  assert (a==1);
}
```

Intuition: the assertion will never fail

Thread A

```c
void foo(void) {
  a = 1;
  b = 1;
}
```

Thread B

```c
void bar(void) {
  while (b == 0){};
  assert (a==1);
}
```

Intuition: the assertion will never fail

⚠ Real execution: given enough tries, the assertion may eventually fail

⤳ in need of defining a *Memory Model*

# Memory Models

Memory interactions behave differently in presence of

- multiple concurrent threads
- data replication in hierarchical and/or distributed memory systems
- deferred communication of updates

Memory Models are a product of negotiating

- restrictions of freedom of implementation to guarantee race related properties
- establishment of freedom of implementation to enable *program* and *machine model* optimizations

$\rightsquigarrow$ Modern Languages include the memory model in their language definition

# Strict Consistency

Motivated by sequential computing, we intuitively implicitly transfer our idea of semantics of memory accesses to concurrent computation. This leads to our idealistic model *Strict Consistency*:

---

**Definition (Strict consistency)**

Independently of which process reads or writes, the value from the most recent write to a location is observable by reads from the respective location immediately *after* the write occurs.

---

Although idealistically desired, practically not existing

⚠ absolute global time problematic

⚠ physically not possible

⤳ strict consistency is too strong to be realistic

# Abandoning absolute time

Thread A

```
void foo(void) {
   a = 1;
   b = 1;
}
```
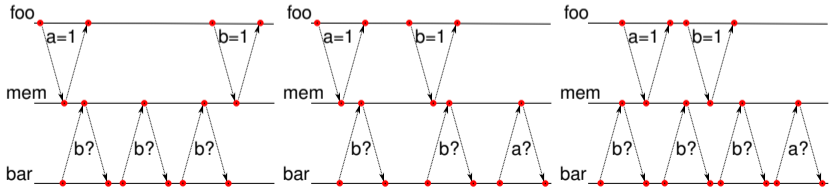
Thread B

```
void bar(void) {
   while (b == 0) {};
   assert(a == 1);
}
```

- initial state of a and b is 0
- A writes a before it writes b
- B should see b go to 1 before executing the assert statement
- the assert statement should always hold
⤳ here correctness means: writing a 1 to a *happens before* reading a 1 in b

Still, *any* of the following may happen:



⤳ Idea: state correctness in terms of what event *may* happen before another one
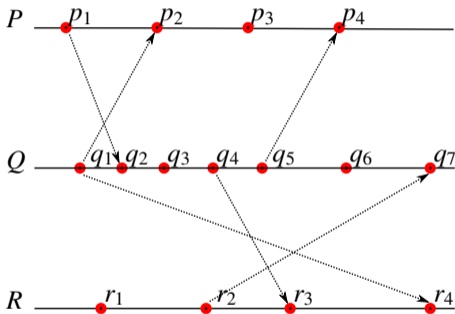
Happend-Before Relation and Diagram

A process as a series of events [Lam78]: Given a distributed system of processes $P, Q, R, \ldots$, each process $P$ consists of events $\bullet p_1, \bullet p_2, \ldots$.

# Events in a Distributed System

A process as a series of events [Lam78]: Given a distributed system of processes $P, Q, R, \ldots$, each process $P$ consists of events $\bullet p_1, \bullet p_2, \ldots$.
Example:



- event $\bullet p_i$ in process $P$ *happened before* $\bullet p_{i+1}$
- if $\bullet p_i$ is an event that sends a message to $Q$ then there is some event $\bullet q_j$ in $Q$ that receives this message and $\bullet p_i$ *happened before* $\bullet q_j$

**Definition**

If an event *p* *happened before* an event *q* then $p \rightarrow q$.

# The Happened-Before Relation

**Definition**

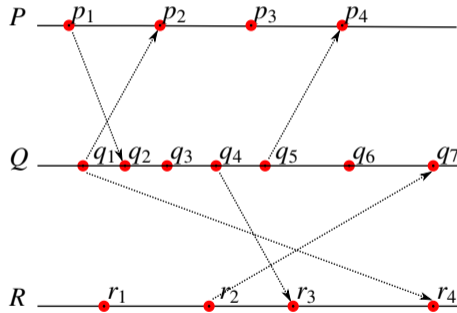If an event $p$ *happened before* an event $q$ then $p \rightarrow q$.

Observe:
- $\rightarrow$ is partial (neither $p \rightarrow q$ or $q \rightarrow p$ may hold)
- $\rightarrow$ is irreflexive ($p \rightarrow p$ never holds)
- $\rightarrow$ is transitive ($p \rightarrow q \wedge q \rightarrow r$ then $p \rightarrow r$)
- $\rightarrow$ is asymmetric (if $p \rightarrow q$ then $\neg(q \rightarrow p)$)

$\rightsquigarrow$ the $\rightarrow$ relation is a *strict partial order*

# Concurrency in Happened-Before Diagrams

Let $a \not\to b$ abbreviate $\neg(a \to b)$.

**Definition**

Two distinct events $p$ and $q$ are said to be *concurrent* if $p \not\to q$ and $q \not\to p$.



- $p_1 \to r_4$ in the example
- $p_3$ and $q_3$ are, in fact, concurrent since $p_3 \not\to q_3$ and $q_3 \not\to p_3$

# **Ordering**

Let $C$ be a *logical clock* i.e. $C$ assigns a *globally unique* time-stamp $C(p)$ to each event $p$.

**Definition (Clock Condition)**

Function $C$ satisfies the *clock condition* **if** for any events $p, q$

$$p \rightarrow q \quad \implies \quad C(p) < C(q)$$

# **Ordering**

Let $C$ be a *logical clock* i.e. $C$ assigns a *globally unique* time-stamp $C(p)$ to each event $p$.

**Definition (Clock Condition)**

Function $C$ satisfies the *clock condition* **if** for any events $p, q$

$$p \rightarrow q \implies C(p) < C(q)$$

For a distributed system the *clock condition* holds iff:

1. $p_i$ and $p_j$ are events of $P$ and $p_i \rightarrow p_j$ then $C(p_i) < C(p_j)$
2. $p$ is the sending of a message by process $P$ and $q$ is the reception of this message by process $Q$ then $C(p) < C(q)$

# Ordering

Let $C$ be a *logical clock* i.e. $C$ assigns a *globally unique* time-stamp $C(p)$ to each event $p$.

**Definition (Clock Condition)**

Function $C$ satisfies the *clock condition* **if** for any events $p, q$

$$p \rightarrow q \quad \implies \quad C(p) < C(q)$$

For a distributed system the *clock condition* holds iff:

1. $p_i$ and $p_j$ are events of $P$ and $p_i \rightarrow p_j$ then $C(p_i) < C(p_j)$
2. $p$ is the sending of a message by process $P$ and $q$ is the reception of this message by process $Q$ then $C(p) < C(q)$
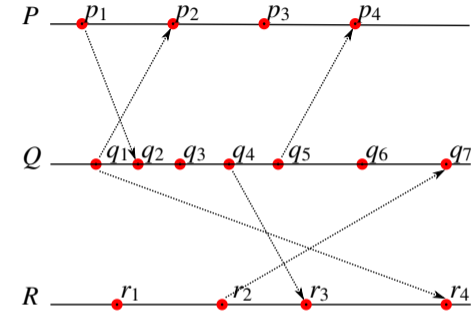
$\rightsquigarrow$ a logical clock $C$ that satisfies the clock condition describes a *total order* $a < b$ (with $C(a) < C(b)$) that *embeds* the strict partial order $\rightarrow$

# Ordering

Let $C$ be a *logical clock* i.e. $C$ assigns a *globally unique* time-stamp $C(p)$ to each event $p$.

## Definition (Clock Condition)

Function $C$ satisfies the *clock condition* **if** for any events $p, q$

$$p \to q \implies C(p) < C(q)$$

For a distributed system the *clock condition* holds iff:

1. $p_i$ and $p_j$ are events of $P$ and $p_i \to p_j$ then $C(p_i) < C(p_j)$
2. $p$ is the sending of a message by process $P$ and $q$ is the reception of this message by process $Q$ then $C(p) < C(q)$

↝ a logical clock $C$ that satisfies the clock condition describes a *total order* $a < b$ (with $C(a) < C(b)$) that *embeds* the strict partial order $\to$

The *set* defined by all $C$ that satisfy the clock condition is exactly the *set* of executions possible in the system.

↝ use the process model and $\to$ to define better consistency model

Given:



| $e$ | | $p_1$ | $p_2$ | $p_3$ | $p_4$ |
|---|---|---|---|---|---|
| $C(e)$ | | | | | |

| $e$ | $q_1$ | $q_2$ | $q_3$ | $q_4$ | $q_5$ | $q_6$ | $q_7$ |
|---|---|---|---|---|---|---|---|
| $C(e)$ | | | | | | | |

| $e$ | | $r_1$ | $r_2$ | $r_3$ | $r_4$ |
|---|---|---|---|---|---|
| $C(e)$ | | | | | |

Given:



| $e$ | $p_1$ | $p_2$ | $p_3$ | $p_4$ |
|---|---|---|---|---|
| $C(e)$ | 1 | 4 | 7 | 12 |

| $e$ | $q_1$ | $q_2$ | $q_3$ | $q_4$ | $q_5$ | $q_6$ | $q_7$ |
|---|---|---|---|---|---|---|---|
| $C(e)$ | 2 | 3 | 5 | 6 | 11 | 13 | 14 |

| $e$ | $r_1$ | $r_2$ | $r_3$ | $r_4$ |
|---|---|---|---|---|
| $C(e)$ | 8 | 9 | 10 | 15 |

# Summing up Happened-Before Relations

We can model concurrency using processes and events:

- there is a *happened-before* relation between the events of each process
- there is a *happened-before* relation between communicating events
- *happened-before* is a strict partial order
- a clock is a total strict order that embeds the *happened-before* partial order

# Memory Consistency Models based on the Happened-Before Relation

# **Happened-Before Based Memory Models**

Idea: use happened-before diagrams to model more relaxed memory models.

Given a path through each of the threads of a program:
- consider the actions of each thread as events of a process
- use more processes to model memory
  - here: one process per variable in memory
- $\leadsto$ concisely represent *some* interleavings

# **Happened-Before Based Memory Models**

Idea: use happened-before diagrams to model more relaxed memory models.

Given a path through each of the threads of a program:
- consider the actions of each thread as events of a process
- use more processes to model memory
  - here: one process per variable in memory
- ⤳ concisely represent *some* interleavings

⤳ We establish a model for *Sequential Consistency*.

# Sequential Consistency

**Definition (Sequential Consistency Condition [Lam78])**

The result of any execution is the same as if the memory operations

- of each individual processor appear in the order specified by its program
- of all processors joined were executed in some sequential order

*Sequential Consistency applied to Multiprocessor Programs:*

Given a program with $n$ threads,

1. for fixed event sequences $p_0^1, p_1^1, \ldots$ and $p_0^2, p_1^2, \ldots$ and $p_0^n, p_1^n, \ldots$ keeping the program order,
2. executions obeying the clock condition on the $p_j^i$,
3. all executions have the same result

Yet, in other words:

- 1 defines the *execution path* of each thread
- each execution mentioned in 2 is one *interleaving* of processes
- 3 declares that the result of running the threads with these interleavings is always the same.

# **Working with Sequential Consistency**

*Sequential Consistency in Multiprocessor Programs:*

Given a program with $n$ threads,

1. for fixed event sequences $p_0^1, p_1^1, \ldots$ and $p_0^2, p_1^2, \ldots$ and $p_0^n, p_1^n, \ldots$ keeping the program order,

2. executions obeying the clock condition on the $p_j^i$,

3. all executions have the same result

Idea for showing that a system is *not* sequentially consistent:

- pick a result obtained from a program run on a SC system
- pick an execution **1** and a total ordering of all operations **2**
- add extra processes to model other system components
- the original order **2** becomes a partial order $\rightarrow$
- show that total orderings $C'$ exist for $\rightarrow$ for which the result differs

## Definition (Sequential Consistency)

1. Memory operations in program order ( $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\mathrm{Op}_i[a] \leq \mathrm{Op}_i[b]' \Rightarrow \mathrm{Op}_i[a] \sqsubseteq \mathrm{Op}_i[b]'$$

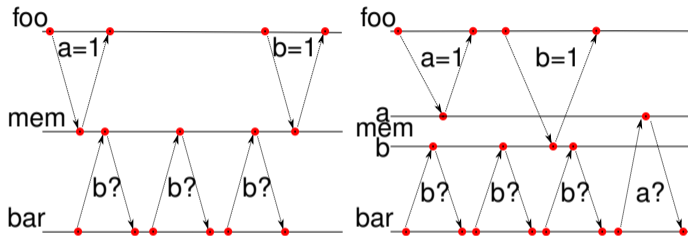2. A load's value is determined by the latest write wrt. memory order

$$val(\mathrm{Ld}_i[a]) = val(\mathrm{St}_j[a] \mid \mathrm{St}_j[a] = \max_{\sqsubseteq} (\{\mathrm{St}_k[a] \mid \mathrm{St}_k[a] \sqsubseteq \mathrm{Ld}_i[a]\}))$$

with

- $\mathrm{Op}_i[a]$ any memory access to address $a$ by CPU $i$
- $\mathrm{Ld}_i[a]$ a load from address $a$ by CPU $i$
- $\mathrm{St}_i[a]$ a store to address $a$ by CPU $i$
- Program order $\leq$ being specified by the control flow of the programs executed by their associated CPUs; only orders operations on the same CPU

## **Weakening the Model**

Observation: more concurrency possible, if we model each memory location separately, i.e. as a different process



Sequential consistency still obeyed:

- the accesses of `foo` to `a` occurs before `b`
- the first two read accesses to `b` are in parallel to `a=1`

Conclusion: There is no observable change if accesses to different memory locations can happen in parallel.

- concisely represent *all* interleavings that are due to variations in timing
- synchronization using time is uncommon for software
- ⇝ a good model for correct behaviors of concurrent programs
- ⇝ program results besides SC results are undesirable (they contain *races*)

# Benefits of Sequential Consistency

- concisely represent *all* interleavings that are due to variations in timing
- synchronization using time is uncommon for software
$\rightsquigarrow$ a good model for correct behaviors of concurrent programs
$\rightsquigarrow$ program results besides SC results are undesirable (they contain *races*)

**Realistic model for simple hardware architectures:**

- sequential consistency model suitable for concurrent processors that acquire *exclusive* access to memory
- processors can speed up computation by using *caches* and still made to maintain sequential consistency

# Benefits of Sequential Consistency

- concisely represent *all* interleavings that are due to variations in timing
- synchronization using time is uncommon for software
- ⤳ a good model for correct behaviors of concurrent programs
- ⤳ program results besides SC results are undesirable (they contain *races*)

**Realistic model for simple hardware architectures:**

- sequential consistency model suitable for concurrent processors that acquire *exclusive* access to memory
- processors can speed up computation by using *caches* and still made to maintain sequential consistency

**Not realistic for elaborate hardware with out-of-order stores:**
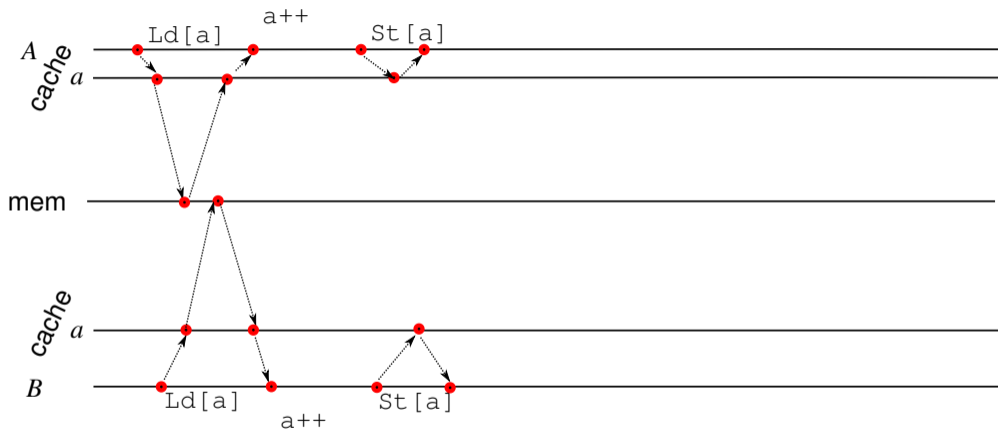
- what other processors see is determined by complex optimizations to cacheline management

⤳ internal workings of caches

Introducing Caches: The MESI Protocol

# **Introducing Caches**

Idea: each cache line one process



Observations:

⚠ naive replication of memory in cache lines creates *incoherency*

**Definition (Cache Coherency)**

1. Memory operations in program order ( $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\text{Op}_i[a] \leq \text{Op}_i[a]' \Rightarrow \text{Op}_i[a] \sqsubseteq \text{Op}_i[a]'$$

2. A load's value is determined by the latest write wrt. memory order

$$val(\text{Ld}_i[a]) = val(\text{St}_j[a] \mid \text{St}_j[a] = \underset{\sqsubseteq}{max} (\{\text{St}_k[a] \mid \text{St}_k[a] \sqsubseteq \text{Ld}_i[a]\}))$$

- This definition superficially looks close to the definition of SC – except that it covers only singular memory locations instead of all memory locations accessed in a program
- Caches and memory can communicate using messaging, following some particular protocol to establish cache coherency
  (⤳ *Cache Coherence Protocol*)

Processors use caches to avoid a costly round-trip to RAM for every memory access.

- programs often access the same memory area repeatedly (e.g. stack)
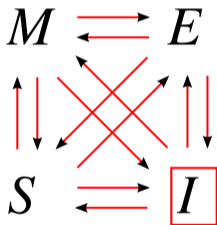- keeping a local mirror image of certain memory regions requires bookkeeping about who has the latest copy



Each cache line is in one of the states $M, E, S, I$:

Processors use caches to avoid a costly round-trip to RAM for every memory access.

- programs often access the same memory area repeatedly (e.g. stack)
- keeping a local mirror image of certain memory regions requires bookkeeping about who has the latest copy



Each cache line is in one of the states $M, E, S, I$:

$I:$ it is *invalid* and is ready for re-use

# The MESI Cache Coherence Protocol: States [PP84]

Processors use caches to avoid a costly round-trip to RAM for every memory access.

- programs often access the same memory area repeatedly (e.g. stack)
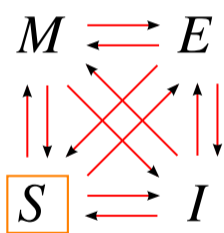- keeping a local mirror image of certain memory regions requires bookkeeping about who has the latest copy



Each cache line is in one of the states $M, E, S, I$:

$I:$ it is *invalid* and is ready for re-use

$S:$ other caches have an identical copy of this cache line, it is *shared*

# The MESI Cache Coherence Protocol: States [PP84]

Processors use caches to avoid a costly round-trip to RAM for every memory access.

- programs often access the same memory area repeatedly (e.g. stack)
- keeping a local mirror image of certain memory regions requires bookkeeping about who has the latest copy
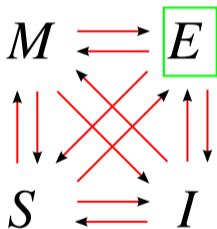


Each cache line is in one of the states $M, E, S, I$:

    *I:* it is *invalid* and is ready for re-use

    *S:* other caches have an identical copy of this cache line, it is *shared*

    *E:* the content is in no other cache; it is *exclusive* to this cache and can be overwritten without consulting other caches

# The MESI Cache Coherence Protocol: States [PP84]

Processors use caches to avoid a costly round-trip to RAM for every memory access.

- programs often access the same memory area repeatedly (e.g. stack)
- keeping a local mirror image of certain memory regions requires bookkeeping about who has the latest copy



Each cache line is in one of the states $M, E, S, I$:

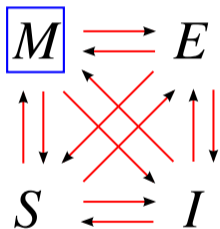$I:$ it is *invalid* and is ready for re-use

$S:$ other caches have an identical copy of this cache line, it is *shared*

$E:$ the content is in no other cache; it is *exclusive* to this cache and can be overwritten without consulting other caches

$M:$ the content is exclusive to this cache and has furthermore been *modified*

Processors use caches to avoid a costly round-trip to RAM for every memory access.

- programs often access the same memory area repeatedly (e.g. stack)
- keeping a local mirror image of certain memory regions requires bookkeeping about who has the latest copy

$$M \rightleftarrows E$$
$$\begin{matrix} \updownarrow & \times & \updownarrow \\ S & \rightleftarrows & I \end{matrix}$$

Each cache line is in one of the states $M, E, S, I$:

- $I$: it is *invalid* and is ready for re-use
- $S$: other caches have an identical copy of this cache line, it is *shared*
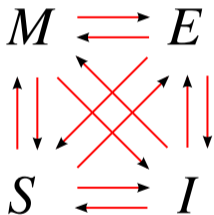- $E$: the content is in no other cache; it is *exclusive* to this cache and can be overwritten without consulting other caches
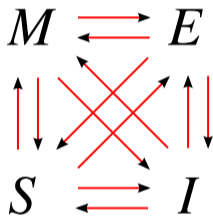- $M$: the content is exclusive to this cache and has furthermore been *modified*

⤳ the global state of cache lines is kept consistent by sending *messages*

# The MESI Cache Coherence Protocol: Messages

Moving data between caches is coordinated by sending messages [McK10]:

- *Read:* sent if CPU needs to read from an address
- *Read Response:* when in state E or S, response to a *Read* message, carries the data for the requested address
- *Invalidate:* asks others to evict a cache line
- *Invalidate Acknowledge:* reply indicating that a cache line has been evicted
- *Read Invalidate:* like *Read* + *Invalidate* (also called "read with intend to modify")
- *Writeback: Read Response* when in state M, as a side effect noticing main memory about modifications to the cacheline, changing sender's state to S

$$M \rightleftarrows E$$
$$S \rightleftarrows I$$

We mostly consider messages between processors. Upon *Read Invalidate*, a processor replies with *Read Response*/*Writeback* before the *Invalidate Acknowledge* is sent.

# MESI Example

Consider how the following code might execute:

| Thread A | |
|---|---|
| `a = 1;` | `// A.1` |
| `b = 1;` | `// A.2` |

| Thread B | |
|---|---|
| `while (b == 0) {};` | `// B.1` |
| `assert(a == 1);` | `// B.2` |

- in all examples, the initial values of variables are assumed to be 0
- suppose that `a` and `b` reside in different cache lines
- assume that a cache line is larger than the variable itself
- we write the content of a cache line as
  - M$x$: modified, with value $x$
  - E$x$: exclusive, with value $x$
  - S$x$: shared, with value $x$
  - I: invalid

# MESI Example (I)

## Thread A

```
a = 1;      // A.1
b = 1;      // A.2
```

## Thread B

```
while (b == 0) {};   // B.1
assert(a == 1);      // B.2
```

| statement | CPU A | | CPU B | | RAM | | message |
|---|---|---|---|---|---|---|---|
| | a | b | a | b | a | b | |
| A.1 | I | I | I | I | 0 | 0 | read invalidate of a from CPU A |
| | I | I | I | I | 0 | 0 | invalidate ack. of a from CPU B |
| | I | I | I | I | 0 | 0 | read response of a=0 from RAM |
| B.1 | M 1 | I | I | I | 0 | 0 | read of b from CPU B |
| | M 1 | I | I | I | 0 | 0 | read response with b=0 from RAM |
| B.1 | M 1 | I | I | E 0 | 0 | 0 | |
| A.2 | M 1 | I | I | E 0 | 0 | 0 | read invalidate of b from CPU A |
| | M 1 | I | I | E 0 | 0 | 0 | read response of b=0 from CPU B |
| | M 1 | S 0 | I | S 0 | 0 | 0 | invalidate ack. of b from CPU B |
| | M 1 | M 1 | I | I | 0 | 0 | |

# MESI Example (II)

## Thread A

```
a = 1;      // A.1
b = 1;      // A.2
```

## Thread B

```
while (b == 0) {};   // B.1
assert(a == 1);      // B.2
```

| statement | CPU A | | CPU B | | RAM | | message |
|---|---|---|---|---|---|---|---|
| | a | b | a | b | a | b | |
| B.1 | M 1 | M 1 | I | I | 0 | 0 | read of b from CPU B |
| | M 1 | M 1 | I | I | 0 | 0 | write back of b=1 from CPU A |
| B.2 | M 1 | S 1 | I | S 1 | 0 | 1 | read of a from CPU B |
| | M 1 | S 1 | I | S 1 | 0 | 1 | write back of a=1 from CPU A |
| | S 1 | S 1 | S 1 | S 1 | 1 | 1 | |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| A.1 | S 1 | S 1 | S 1 | S 1 | 1 | 1 | invalidate of a from CPU A |
| | S 1 | S 1 | I | S 1 | 1 | 1 | invalidate ack. of a from CPU B |
| | M 1 | S 1 | I | S 1 | 1 | 1 | |

# MESI Example: Happened Before Model

Idea: each cache line one process, A caches b=0 as E, B caches a=0 as E
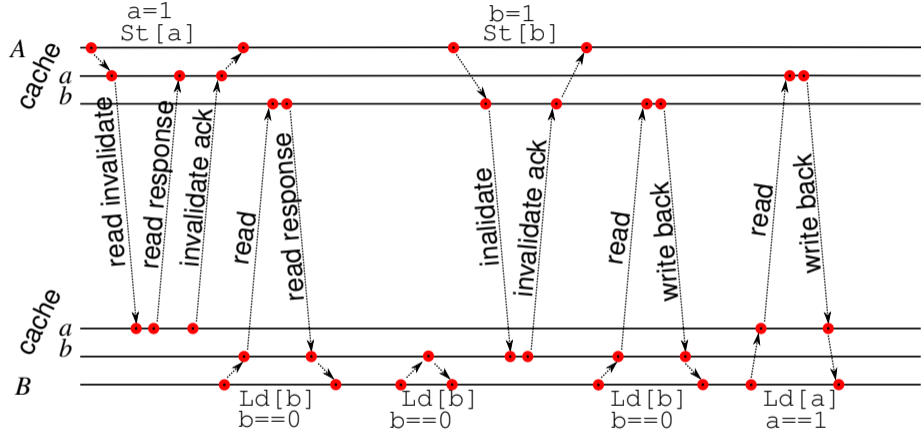


Observations:

- each memory access must complete before executing next instruction ⤳ add edge

## MESI Example: Happened Before Model

Idea: each cache line one process, A caches b=0 as E, B caches a=0 as E



Observations:

- each memory access must complete before executing next instruction ⤳ add edge
- second execution of test b==0 stays within cache ⤳ no traffic

*Sequential Consistency:*

- specifies that the system must appear to execute all threads' loads and stores to *all memory locations* in a total order that respects the program order of each thread
- a characterization of well-behaved programs
- a model for differing speed of execution
- for fixed paths through the threads *and* a total order between accesses to the same variables: executions can be illustrated by a happened-before diagram with one process per variable

*Cache Coherency:*

- A *cache coherent* system must appear to execute all threads' loads and stores to a *single memory location* in a total order that respects the program order of each thread
- MESI cache coherence protocol ensures SC for processors with caches

Introducing Store Buffers: Out-Of-Order Stores
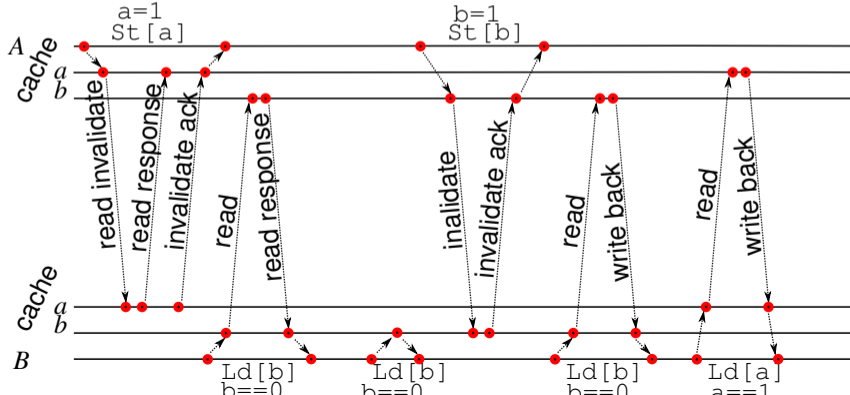
# Out-of-Order Execution

⚠ performance problem: writes always stall

| Thread A | Thread B |
|---|---|
| `a = 1;    // A.1`<br>`b = 1;    // A.2` | `while (b == 0) {};  // B.1`<br>`assert(a == 1);     // B.2` |

# Out-of-Order Execution

⚠ performance problem: writes always stall
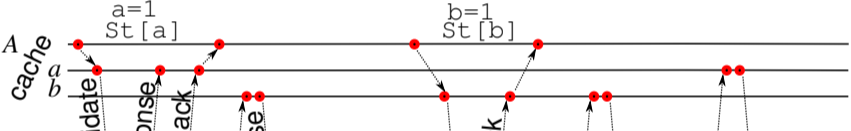
| Thread A | Thread B |
|---|---|
| ```
a = 1;    // A.1
b = 1;    // A.2
``` | ```
while (b == 0) {};  // B.1
assert(a == 1);     // B.2
``` |



↝ CPU **A** should continue executing after `a = 1;`

# Store Buffers

⚠ *Abstract Machine Model:* defines semantics of memory accesses



- put *each* store into a *store buffer* and continue execution
- Store buffers apply stores in various orders:
  - ▶ FIFO (Sparc/x86-*TSO*)
  - ▶ unordered (Sparc *PSO*)
- ⚠ program order still needs to be observed locally
  - ▶ store buffer snoops read channel and
  - ▶ on matching address, returns the youngest value in buffer

## Definition (Total Store Order)

**1** The store order wrt. memory ( $\sqsubseteq$ ) is total

$$\forall_{a,b \,\in\, addr \; i,j \,\in\, CPU} \quad (\mathtt{St}_i[a] \sqsubseteq \mathtt{St}_j[b]) \vee (\mathtt{St}_j[b] \sqsubseteq \mathtt{St}_i[a])$$

**2** Stores in program order ( $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\mathtt{St}_i[a] \leq \mathtt{St}_i[b] \Rightarrow \mathtt{St}_i[a] \sqsubseteq \mathtt{St}_i[b]$$

**3** Loads preceding an other operation (wrt. program order $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\mathtt{Ld}_i[a] \leq \mathtt{Op}_i[b] \Rightarrow \mathtt{Ld}_i[a] \sqsubseteq \mathtt{Op}_i[b]$$

**4** A load's value is determined by the latest write as observed by the local CPU

$$val(\mathtt{Ld}_i[a]) = val(\mathtt{St}_j[a] \mid \mathtt{St}_j[a] = \max_{\sqsubseteq} (\{\mathtt{St}_k[a] \mid \mathtt{St}_k[a] \sqsubseteq \mathtt{Ld}_i[a]\} \cup \{\mathtt{St}_i[a] \mid \mathtt{St}_i[a] \leq \mathtt{Ld}_i[a]\}))$$

Particularly, one ordering property from SC is not guaranteed:

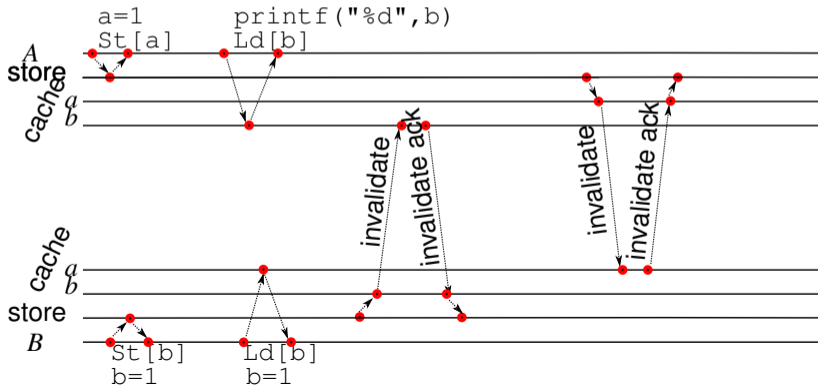$$\mathtt{St}_i[a] \leq \mathtt{Ld}_i[b] \not\Rightarrow \mathtt{St}_i[a] \sqsubseteq \mathtt{Ld}_i[b]$$

⚠ Local stores may be observed earlier by local loads then from somewhere else!

# Happened-Before Model for TSO

| Thread A | Thread B |
|---|---|
| ```
a = 1;
printf("%d",b);
``` | ```
b = 1;
printf("%d",a);
``` |

Assume cache A contains: a: S0, b: S0, cache B contains: a: S0, b: S0

## TSO in the Wild: x86

The x86 CPU, powering desktops and servers around the world is a common representative of a TSO Memory Model based CPU.

- FIFO store buffers keep quite strong consistency properties
- The major obstacle to Sequential Consistency is

$$\text{St}_i[a] \leq \text{Ld}_i[b] \quad \not\Rightarrow \quad \text{St}_i[a] \sqsubseteq \text{Ld}_i[b]$$

  ▶ modern x86 CPUs provide the `mfence` instruction
  ▶ `mfence` orders all memory instructions:

  $$\text{Op}_i \leq \textit{mfence}() \leq \text{Op}_i{}' \quad \Rightarrow \quad \text{Op}_i \sqsubseteq \text{Op}_i{}'$$

- a fence between write and loads gives sequentially consistent CPU behavior (and is as slow as a CPU without store buffer)
$\rightsquigarrow$ use fences only when necessary

## Definition (Partial Store Order)

1. The store order wrt. memory ( $\sqsubseteq$ ) is total

$$\forall_{a,b \in addr\ i,j \in CPU} \quad (\text{St}_i[a] \sqsubseteq \text{St}_j[b]) \vee (\text{St}_j[b] \sqsubseteq \text{St}_i[a])$$

2. Fenced stores in program order ( $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\text{St}_i[a] \leq \text{sfence}() \leq \text{St}_i[b] \Rightarrow \text{St}_i[a] \sqsubseteq \text{St}_i[b]$$

3. Stores to the same address in program order ( $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\text{St}_i[a] \leq \text{St}_i[a]' \Rightarrow \text{St}_i[a] \sqsubseteq \text{St}_i[a]'$$

4. Loads preceding another operation (wrt. program order $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\text{Ld}_i[a] \leq \text{Op}_i[b] \Rightarrow \text{Ld}_i[a] \sqsubseteq \text{Op}_i[b]$$

5. A load's value is determined by the latest write as observed by the local CPU

$$val(\text{Ld}_i[a]) = val(\text{St}_j[a] \mid \text{St}_j[a] = \max_{\sqsubseteq} (\{\text{St}_k[a] \mid \text{St}_k[a] \sqsubseteq \text{Ld}_i[a]\} \cup \{\text{St}_i[a] \mid \text{St}_i[a] \leq \text{Ld}_i[a]\}))$$

⚠ Now also stores are not guaranteed to be in order any more:

$$\text{St}_i[a] \leq \text{St}_i[b] \not\Rightarrow \text{St}_i[a] \sqsubseteq \text{St}_i[b]$$

⤳ What about sequential consistency for the whole system?

## Happened-Before Model for PSO

| Thread A |
|----------|
| ```
a = 1;
b = 1;
``` |

| Thread B |
|----------|
| ```
while (b == 0) {};
assert(a == 1);
``` |

Assume cache A contains: a: S0, b: E0, cache B contains: a: S0, b: I

Overtaking of messages *may be desirable* and does not need to be prohibited in general.

- generalized store buffers render programs incorrect that assume sequential consistency between *different* CPUs
- whenever a store in front of another operation in one CPU must be observable in this order *by a different CPU*, an explicit *write barrier* has to be inserted
  - a write barrier marks all current store operations in the store buffer
  - the next store operation is only executed when all marked stores in the buffer have completed

| Thread A |
|---|
| ```
a = 1;
sfence();
b = 1;
``` |

| Thread B |
|---|
| ```
while (b == 0) {};
assert(a == 1);
``` |

Assume cache A contains: a: S0, b: E0, cache B contains: a: S0, b: I
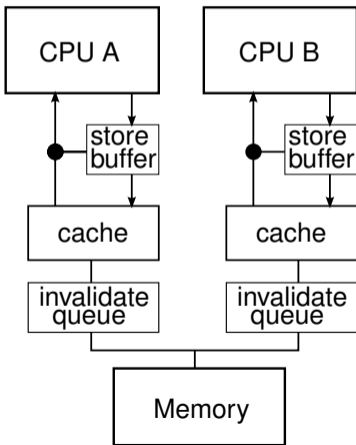
Further weakening the model: O-o-O Reads

## Relaxed Memory Order

Communication of cache updates is still costly:

- a cache-intense computation can fill up store buffers in CPUs
- ↝ waiting for invalidation acknoledgements may still happen
- invalidation acknoledgements are delayed on busy caches

CPU A — store buffer — cache — invalidate queue

CPU B — store buffer — cache — invalidate queue

Memory

- ↝ immediately acknowledge an invalidation and apply it later
  - put each invalidate message into an *invalidate queue*
  - if a *MESI message* needs to be sent regarding a cache line in the invalidate queue then wait until the line is invalidated
- ⚠ local loads and stores do *not* consult the invalidate queue
- ↝ What about sequential consistency?

**Definition (Relaxed Memory Order)**

**1** Fenced memory accesses in program order ( $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\text{Op}_i[a] \leq \text{mfence()} \leq \text{Op}_i[b] \Rightarrow \text{Op}_i[a] \sqsubseteq \text{Op}_i[b]$$

**2** Stores to the same address in program order ( $\leq$ ) are embedded into the memory order ( $\sqsubseteq$ )

$$\text{Op}_i[a] \leq \text{St}_i[a]' \Rightarrow \text{Op}_i[a] \sqsubseteq \text{St}_i[a]'$$

**3** Operations dependent on a load (wrt. *dependence* $\rightarrow$ ) are embedded in the memory order ( $\sqsubseteq$ )

$$\text{Ld}_i[a] \rightarrow \text{Op}_i[b] \Rightarrow \text{Ld}_i[a] \sqsubseteq \text{Op}_i[b]$$

**4** A load's value is determined by the latest write as observed by the local CPU

$$val(\text{Ld}_i[a]) = val(\text{St}_j[a] \mid \text{St}_j[a] = \max_{\sqsubseteq} (\{\text{St}_k[a] \mid \text{St}_k[a] \sqsubseteq \text{Ld}_i[a]\} \cup \{\text{St}_i[a] \mid \text{St}_i[a] \leq \text{Ld}_i[a]\}))$$

⚠ Now we need the notion of *dependence* $\rightarrow$ :

- Memory access to the same address: $\quad \text{St}_i[a] \leq \text{Ld}_i[a] \quad \Rightarrow \quad \text{St}_i[a] \rightarrow \text{Ld}_i[a]$
- Register reads are dependent on latest register writes:

$$\text{Ld}_i[a]'' = \max_{\leq} (\text{Ld}_i[a]' \mid targetreg(\text{Ld}_i[a]') = srcreg(\text{St}_i[b]) \land \text{Ld}_i[a]' \leq \text{St}_i[b]) \quad \Rightarrow \quad \text{Ld}_i[a]'' \rightarrow \text{St}_i[b]$$

- Stores within branched blocks are dependent on branch conditionals:

$$(\text{Op}_i[a] \leq \text{St}_i[b]) \land \text{Op}_i[a] \rightarrow condbranch \leq \text{St}_i[b] \quad \Rightarrow \quad \text{Op}_i[a] \rightarrow \text{St}_i[b]$$
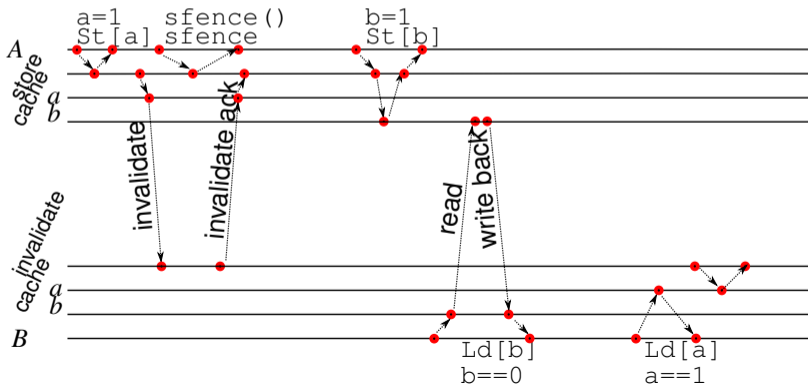
# Happened-Before Model for Invalidate Queues

| Thread A |
|---|
| ```
a = 1;
sfence();
b = 1;
``` |

| Thread B |
|---|
| ```
while (b == 0) {};
assert(a == 1);
``` |

Assume cache A contains: a: S0, b: E0, cache B contains: a: S0, b: I

Read accesses do not consult the invalidate queue.

- might read an out-of-date value
- need a way to establish sequential consistency between writes of other processors and local reads
- insert an explicit *read barrier* before the read access
  - a read barrier marks all entries in the invalidate queue
  - the next read operation is only executed once all marked invalidations have completed
- a read barrier *before* each read gives sequentially consistent read behavior (and is as slow as a system without invalidate queue)

$\rightsquigarrow$ match each write barrier in one process with a read barrier in another process
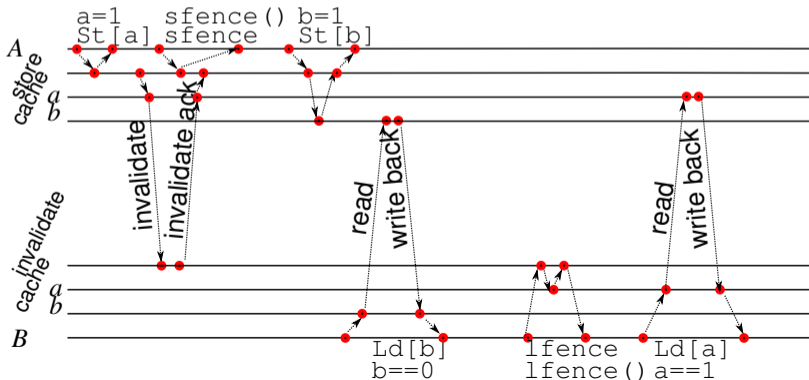
# Happened-Before Model for Read Barriers



| Thread A | Thread B |
|---|---|
| ```a = 1;``` ```sfence();``` ```b = 1;``` | ```while (b == 0) {};``` ```lfence();``` ```assert(a == 1);``` |

Example: The Dekker Algorithm on RMO Systems

## Using Memory Barriers: the Dekker Algorithm

Mutual exclusion of *two* processes with busy waiting.

```
//flag[] is boolean array; and turn is an integer
flag[0] = false;
flag[1] = false;
turn    = 0;   // or 1
```

```
P0:
flag[0] = true;
while (flag[1] == true)
  if (turn != 0) {
     flag[0] = false;
     while (turn != 0) {
       // busy wait
     }
     flag[0] = true;
  }
// critical section
turn    = 1;
flag[0] = false;
```

## Using Memory Barriers: the Dekker Algorithm

Mutual exclusion of *two* processes with busy waiting.

```
//flag[] is boolean array; and turn is an integer
flag[0] = false;
flag[1] = false;
turn    = 0;    // or 1
```

```
P0:
flag[0] = true;
while (flag[1] == true)
  if (turn != 0) {
     flag[0] = false;
     while (turn != 0) {
       // busy wait
     }
     flag[0] = true;
  }
// critical section
turn    = 1;
flag[0] = false;
```

```
P1:
flag[1] = true;
while (flag[0] == true)
  if (turn != 1) {
     flag[1] = false;
     while (turn != 1) {
       // busy wait
     }
     flag[1] = true;
  }
// critical section
turn    = 0;
flag[1] = false;
```

## The Idea Behind Dekker

Communication via three variables:

- flag[i]==true process $P_i$ wants to enter its critical section
- turn==i process $P_i$ has priority when both want to enter

```
P0:
flag[0] = true;
while (flag[1] == true)
  if (turn != 0) {
      flag[0] = false;
      while (turn != 0) {
        // busy wait
      }
      flag[0] = true;
  }
// critical section
turn    = 1;
flag[0] = false;
```

In process $P_i$:

- if $P_{1-i}$ does not want to enter, proceed immediately to the critical section

## **The Idea Behind Dekker**

Communication via three variables:

- $flag[i]==true$ process $P_i$ wants to enter its critical section
- $turn==i$ process $P_i$ has priority when both want to enter

```
P0:
flag[0] = true;
while (flag[1] == true)
  if (turn != 0) {
      flag[0] = false;
      while (turn != 0) {
        // busy wait
      }
      flag[0] = true;
  }
// critical section
turn    = 1;
flag[0] = false;
```

In process $P_i$:

- if $P_{1-i}$ does not want to enter, proceed immediately to the critical section
- $\rightsquigarrow$ flag[i] is a *lock* and may be implemented as such

## The Idea Behind Dekker

Communication via three variables:

- flag[i]==true process $P_i$ wants to enter its critical section
- turn==i process $P_i$ has priority when both want to enter

```
P0:
flag[0] = true;
while (flag[1] == true)
  if (turn != 0) {
      flag[0] = false;
      while (turn != 0) {
        // busy wait
      }
      flag[0] = true;
  }
// critical section
turn    = 1;
flag[0] = false;
```

In process $P_i$:

- if $P_{1-i}$ does not want to enter, proceed immediately to the critical section
- ⇝ flag[i] is a *lock* and may be implemented as such
- if $P_{1-i}$ also wants to enter, wait for turn to be set to i

## The Idea Behind Dekker

Communication via three variables:

- flag[i]==true process $P_i$ wants to enter its critical section
- turn==i process $P_i$ has priority when both want to enter

```
P0:
flag[0] = true;
while (flag[1] == true)
  if (turn != 0) {
      flag[0] = false;
      while (turn != 0) {
        // busy wait
      }
      flag[0] = true;
  }
// critical section
turn    = 1;
flag[0] = false;
```

In process $P_i$:

- if $P_{1-i}$ does not want to enter, proceed immediately to the critical section
- ⇝ flag[i] is a *lock* and may be implemented as such
  - if $P_{1-i}$ also wants to enter, wait for turn to be set to i
  - while waiting for turn, reset flag[i] to enable $P_{1-i}$ to progress

Problem: Dekker's algorithm requires sequential consistency.
Idea: insert memory barriers between all variables common to both threads.

## Dekker's Algorithm and RMO

Problem: Dekker's algorithm requires sequential consistency.
Idea: insert memory barriers between all variables common to both threads.

```
P0:
flag[0] = true;
sfence();
while (lfence(), flag[1] == true)
  if (lfence(), turn != 0) {
     flag[0] = false;
     sfence();
     while (lfence(), turn != 0){
       // busy wait
     }
     flag[0] = true;
     sfence();
  }
// critical section
turn    = 1;
sfence();
flag[0] = false; sfence();
```

- insert a load memory barrier lfence()
  in front of every read from common
  variables

## **Dekker's Algorithm and RMO**

Problem: Dekker's algorithm requires sequential consistency.
Idea: insert memory barriers between all variables common to both threads.

```
P0:
flag[0] = true;
sfence();
while (lfence(), flag[1] == true)
  if (lfence(), turn != 0) {
     flag[0] = false;
     sfence();
     while (lfence(), turn != 0){
       // busy wait
     }
     flag[0] = true;
     sfence();
  }
// critical section
turn    = 1;
sfence();
flag[0] = false; sfence();
```

- insert a load memory barrier `lfence()` in front of every read from common variables
- insert a write memory barrier `sfence()` after writing a variable that is read in the other thread

## Dekker's Algorithm and RMO

Problem: Dekker's algorithm requires sequential consistency.
Idea: insert memory barriers between all variables common to both threads.

```
P0:
flag[0] = true;
sfence();
while (lfence(), flag[1] == true)
  if (lfence(), turn != 0) {
     flag[0] = false;
     sfence();
     while (lfence(), turn != 0){
       // busy wait
     }
     flag[0] = true;
     sfence();
  }
// critical section
turn    = 1;
sfence();
flag[0] = false; sfence();
```

- insert a load memory barrier `lfence()` in front of every read from common variables
- insert a write memory barrier `sfence()` after writing a variable that is read in the other thread
- the `lfence()` of the first iteration of each loop may be combined with the preceding `sfence()` to an `mfence()`

## Summary: Relaxed Memory Models

Highly optimized CPUs may use a *relaxed memory model*:

- reads and writes are not synchronized unless requested by the user
- many kinds of memory barriers exist with subtle differences

$\rightsquigarrow$ ARM, PowerPC, Alpha, ia-64, even x86 ($\rightsquigarrow$ SSE Write Combining)

$\rightsquigarrow$ memory barriers are the "lowest-level" of synchronization

## **Discussion**

Memory barriers reside at the lowest level of synchronization primitives.

## Discussion

Memory barriers reside at the lowest level of synchronization primitives.

Where are they useful?

- when blocking should not de-schedule threads
- when several processes implement automata and coordinate their transitions via common synchronized variables
- ⇝ protocol implementations
- ⇝ OS provides synchronization facilities based on memory barriers

Why might they not be appropriate?

- difficult to get right, best suited for specific well-understood algorithms
- often synchronization with locks is as fast and easier
- too many fences are costly if store/invalidate buffers are bottleneck

Before Optimization

```
int x = 0;
for (int i=0;i<100;i++){
    x = 1;
    printf("%d",x);
}
```

# Memory Models and Compilers

Before Optimization

```
int x = 0;
for (int i=0;i<100;i++){
    x = 1;
    printf("%d",x);
}
```

After Optimization

```
int x = 1;
for (int i=0;i<100;i++){
    printf("%d",x);
}
```

**Standard Program Optimizations**

comprises *loop-invariant code motion* and *dead store elimination*, e.g.

# Memory Models and Compilers

Before Optimization

```c
int x = 0;
for (int i=0; i<100; i++) {
    x = 1;
    printf("%d", x);
}
```

After Optimization

```c
int x = 1;
for (int i=0; i<100; i++) {
    printf("%d", x);
}
```

## Standard Program Optimizations

comprises *loop-invariant code motion* and *dead store elimination*, e.g.

⚠ having another thread executing `x = 0;` changes observable behaviour depending on optimizing or not

⤳ Compiler also depends on consistency guarantees
⤳ Demand for Memory Models on language level

Keeping semantics I

```c
int x = 0;
for (int i=0;i<100;i++){
    sfence();
    x = 1;
    printf("%d",x);
}
```

# Memory Models and C-Compilers

Keeping semantics I

```
int x = 0;
for (int i=0;i<100;i++){
    sfence();
    x = 1;
    printf("%d",x);
}
```

Keeping semantics II

```
volatile int x = 0;
for (int i=0;i<100;i++){
    x = 1;
    printf("%d",x);
}
```

- Compilers may also reorder store instructions
- Write barriers keep the compiler from reordering across
- The specification of `volatile` keeps the *C-Compiler* from reordering memory accesses to this address

# Memory Models and C-Compilers

Keeping semantics I

```
int x = 0;
for (int i=0;i<100;i++){
    sfence();
    x = 1;
    printf("%d",x);
}
```

Keeping semantics II

```
volatile int x = 0;
for (int i=0;i<100;i++){
    x = 1;
    printf("%d",x);
}
```

- Compilers may also reorder store instructions
- Write barriers keep the compiler from reordering across
- The specification of `volatile` keeps the *C-Compiler* from reordering memory accesses to this address
- *Java*-Compilers even generate barriers around accesses to `volatile` variables

**Learning Outcomes**

1. Strict Consistency
2. Happened-before Relation
3. Sequential Consistency
4. The MESI Cache Model
5. TSO: FIFO store buffers
6. PSO: store buffers
7. RMO: invalidate queues
8. Reestablishing Sequential Consistency with memory barriers
9. Dekker's Algorithm for Mutual Exclusion

# Future Many-Core Systems: NUMA

## Many-Core Machines' Read Responses congest the bus

In that case: Intel's *MESI*F (Forward) to reduce communication overhead.

⚠ But in general, Symmetric multi-processing (SMP) has its limits:
- a memory-intensive computation may cause contention on the bus
- the speed of the bus is limited since the electrical signal has to travel to all participants
- point-to-point connections are faster than a bus, but do not provide possibility of forming consensus

# Future Many-Core Systems: NUMA

**Many-Core Machines' Read Responses congest the bus**

In that case: Intel's *MESIF* (Forward) to reduce communication overhead.

⚠ But in general, Symmetric multi-processing (SMP) has its limits:
- a memory-intensive computation may cause contention on the bus
- the speed of the bus is limited since the electrical signal has to travel to all participants
- point-to-point connections are faster than a bus, but do not provide possibility of forming consensus

⤳ use a bus locally, use point-to-point links globally: *NUMA*

# Future Many-Core Systems: NUMA

## Many-Core Machines' Read Responses congest the bus

In that case: Intel's *MESI*F (Forward) to reduce communication overhead.

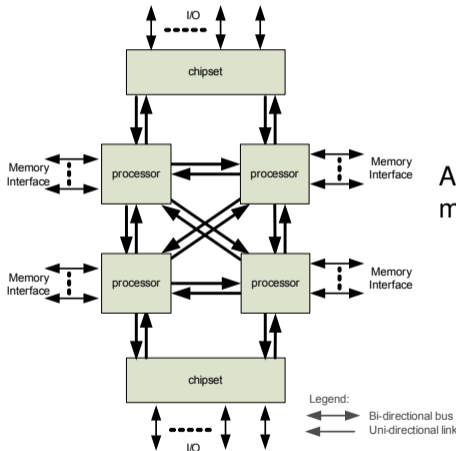⚠ But in general, Symmetric multi-processing (SMP) has its limits:
- a memory-intensive computation may cause contention on the bus
- the speed of the bus is limited since the electrical signal has to travel to all participants
- point-to-point connections are faster than a bus, but do not provide possibility of forming consensus

⇝ use a bus locally, use point-to-point links globally: *NUMA*
- *non-uniform memory access* partitions the memory amongst CPUs
- a directory states which CPU holds a memory region
- Interprocess communication between Cache-Controllers (*ccNUMA*): onchip on Opteron or in chipset on Itanium

## Overhead of NUMA Systems

Communication overhead in a NUMA system.



source: [Int09]

- Processors in a NUMA system may be fully or partially connected.
- The directory of who stores an address is partitioned amongst processors.

A cache miss that cannot be satisfied by the local memory at $A$:

- $A$ sends a retrieve request to processor $B$ owning the directory
- $B$ tells the processor $C$ who holds the content
- $C$ sends data (or status) to $A$ and sends acknowledge to $B$
- $B$ completes transmission by an acknowledge to $A$

# References

Intel.
An introduction to the intel quickpath interconnect.
Technical Report 320412, 2009.

Leslie Lamport.
Time, Clocks, and the Ordering of Events in a Distributed System.
*Commun. ACM*, 21(7):558–565, July 1978.

Paul E. McKenny.
Memory Barriers: a Hardware View for Software Hackers.
Technical report, Linux Technology Center, IBM Beaverton, June 2010.

Mark S. Papamarcos and Janak H. Patel.
A low overhead coherence solution for multiprocessors with private cache memories.
In *In Proc. 11th ISCA*, pages 348–354, 1984.

Daniel J. Sorin, Mark D. Hill, and David A. Wood.
*A Primer on Memory Consistency and Cache Coherence.*
Morgan & Claypool Publishers, 1st edition, 2011.

CORPORATE SPARC International, Inc.
*The SPARC Architecture Manual: Version 8.*
Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1992.

CORPORATE SPARC International, Inc.
*The SPARC Architecture Manual (Version 9).*
Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1994.

- *Sequential Consistency* specifies that the system must appear to execute all threads' loads and stores to *all memory locations* in a total order that respects the program order of each thread
- A *cache coherent* system must appear to execute all threads' loads and stores to a *single memory location* in a total order that respects the program order of each thread

All discussed memory models (SC, TSO, PSO, RMO) provide cache coherence!

# **Programming Languages**

Concurrency: Atomic Executions, Locks and Monitors

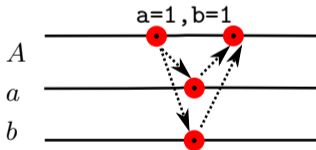Dr. Michael Petter
Winter 2019

# Why Memory Barriers are not Enough

Often, *multiple memory locations* may only be modified exclusively by one thread during a computation.

- use barriers to implement automata that ensure *mutual exclusion*
- ⤳ generalize the re-occurring *concept* of enforcing mutual exclusion

# **Why Memory Barriers are not Enough**

Often, *multiple memory locations* may only be modified exclusively by one thread during a computation.

- use barriers to implement automata that ensure *mutual exclusion*

⤳ generalize the re-occurring *concept* of enforcing mutual exclusion

Needed: interaction with *multiple memory locations* within a *single step*:

## Atomic Executions

A concurrent program consists of several threads that share *resources*:

- resources can be *memory locations* or *memory mapped I/O*
  - ▶ a file can be modified through a shared handle, e.g.
- usually *invariants* must be retained wrt. resources
  - ▶ e.g. a head and tail pointer must delimit a linked list
  - ▶ an invariant may span *multiple* resources
  - ▶ during an update, the invariant may be temporarily *locally broken*
- ↝ multiple resources must be updated together to ensure the invariant

# Atomic Executions

A concurrent program consists of several threads that share *resources*:

- resources can be *memory locations* or *memory mapped I/O*
  - ▶ a file can be modified through a shared handle, e.g.
- usually *invariants* must be retained wrt. resources
  - ▶ e.g. a head and tail pointer must delimit a linked list
  - ▶ an invariant may span *multiple* resources
  - ▶ during an update, the invariant may be temporarily *locally broken*
⤳ multiple resources must be updated together to ensure the invariant

Ideally, a sequence of operations that update shared resources should be *atomic* [Harris et al.(2010)Harris, Larus, and Rajwar]. This would ensure that the invariant never seems to be broken.
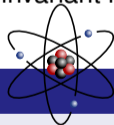
# Atomic Executions

A concurrent program consists of several threads that share *resources*:
- resources can be *memory locations* or *memory mapped I/O*
  - ▶ a file can be modified through a shared handle, e.g.
- usually *invariants* must be retained wrt. resources
  - ▶ e.g. a head and tail pointer must delimit a linked list
  - ▶ an invariant may span *multiple* resources
  - ▶ during an update, the invariant may be temporarily *locally broken*
- ⤳ multiple resources must be updated together to ensure the invariant

Ideally, a sequence of operations that update shared resources should be *atomic* [Harris et al.(2010)Harris, Larus, and Rajwar]. This would ensure that the invariant never seems to be broken.

**Definition (Atomic Execution)**

A computation forms an *atomic execution* if its effect can only be *observed* as a single transformation on the memory.

# Overview

We will address the *established* ways of managing synchronization. The presented techniques

- are available on most platforms
- likely to be found in most existing (concurrent) software
- provide solutions to common concurrency tasks
- are the source of common concurrency problems

The techniques are applicable to C, C++ (pthread), Java, C# and other imperative languages.

## Overview

We will address the *established* ways of managing synchronization. The presented techniques

- are available on most platforms
- likely to be found in most existing (concurrent) software
- provide solutions to common concurrency tasks
- are the source of common concurrency problems

The techniques are applicable to C, C++ (pthread), Java, C# and other imperative languages.

### Learning Outcomes

1. Principle of Atomic Executions
2. Wait-Free Algorithms based on Atomic Operations
3. Locks: Mutex, Semaphore, and Monitor
4. Deadlocks: Concept and Prevention

**Wait-Free Atomic Executions**

Which operations on a CPU are atomic? (`j`,`k` and `tmp` are registers)

**Program 1**

```
i++;
```

**Program 2**

```
j = i;
i = i+k;
```

**Program 3**

```
int tmp = i;
i = j;
j = tmp;
```

Which operations on a CPU are atomic? (`j,k` and `tmp` are registers)

**Program 1**

```
i++;
```

**Program 2**

```
j = i;
i = i+k;
```

**Program 3**

```
int tmp = i;
i = j;
j = tmp;
```

Answer:

- none by default (even without store and invalidate buffers, *why?* )

# Wait-Free Updates

Which operations on a CPU are atomic? (`j,k` and `tmp` are registers)

### Program 1

```
i++;
```

### Program 2

```
j = i;
i = i+k;
```

### Program 3

```
int tmp = i;
i = j;
j = tmp;
```

Answer:

- none by default (even without store and invalidate buffers, *why?* )
- ⚠ The load and store (even `i++`'s) may be interleaved with a store from another processor.

# Wait-Free Updates

Which operations on a CPU are atomic? (`j`,`k` and `tmp` are registers)

| **Program 1** |
| --- |
| `i++;` |

| **Program 2** |
| --- |
| `j = i;`<br>`i = i+k;` |

| **Program 3** |
| --- |
| `int tmp = i;`<br>`i = j;`<br>`j = tmp;` |

Answer:

- none by default (even without store and invalidate buffers, *why?* )
- ⚠ The load and store (even `i++`'s) may be interleaved with a store from another processor.

All of the programs *can* be made atomic executions (e.g. on x86):

- `i` must be in memory
- Idea: *lock the cache bus* for an address for the duration of an instruction

# Wait-Free Updates

Which operations on a CPU are atomic? (`j`,`k` and `tmp` are registers)

**Program 1**
```
i++;
```

**Program 2**
```
j = i;
i = i+k;
```

**Program 3**
```
int tmp = i;
i = j;
j = tmp;
```

Answer:

- none by default (even without store and invalidate buffers, *why?* )
- ⚠ The load and store (even `i++`'s) may be interleaved with a store from another processor.

All of the programs *can* be made atomic executions (e.g. on x86):

- `i` must be in memory
- Idea: *lock the cache bus* for an address for the duration of an instruction

**Program 1**
```
lock inc [addr_i]
```

**Program 2 (fetch-and-add)**
```
mov eax,reg_k
lock xadd [addr_i],eax
mov reg_j,eax
```

**Program 3 (atomic-exchange)**
```
lock xchg [addr_i],reg_j
```

# Wait-Free Bumper-Pointer Allocation

Garbage collectors often use a *bumper pointer* to allocated memory:

## Bumper Pointer Allocation

```
char heap[1<<20];
char* firstFree = &heap[0];

char* alloc(int size) {
  char* start = firstFree;
  firstFree = firstFree + size;

  if (start+size>sizeof(heap)) garbage_collect();
  return start;
}
```

- `firstFree` points to the first unused byte
- each allocation reserves the next `size` bytes in `heap`

# Wait-Free Bumper-Pointer Allocation

Garbage collectors often use a *bumper pointer* to allocated memory:

**Bumper Pointer Allocation**

```c
char heap[1<<20];
char* firstFree = &heap[0];

char* alloc(int size) {
  char* start;
  asm("lock; xadd %0, %1" :"=r"(start),"=m"(firstFree):
      "0"(size),"m"(firstFree) :"memory");
  if (start+size>sizeof(heap)) garbage_collect();
  return start;
}
```

- `firstFree` points to the first unused byte
- each allocation reserves the next `size` bytes in `heap`

Thread-safe implementation:

- `alloc`'s core functionality matches Program 2: fetch-and-add
- inline assembler (GCC/AT&T syntax in the example)

# Marking Statements as Atomic

Rather than writing assembler: use *made-up* keyword `atomic`:

## Program 1

```
atomic {
  i++;
}
```

## Program 2

```
atomic {
  j = i;
  i = i+k;
}
```

## Program 3

```
atomic {
  int tmp = i;
  i = j;
  j = tmp;
}
```

# Marking Statements as Atomic

Rather than writing assembler: use *made-up* keyword `atomic`:

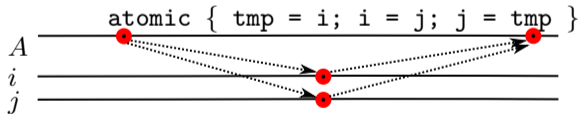**Program 1**
```
atomic {
  i++;
}
```

**Program 2**
```
atomic {
  j = i;
  i = i+k;
}
```

**Program 3**
```
atomic {
  int tmp = i;
  i = j;
  j = tmp;
}
```

The statements in an `atomic` block execute as *atomic execution*:

# Marking Statements as Atomic

Rather than writing assembler: use *made-up* keyword `atomic`:
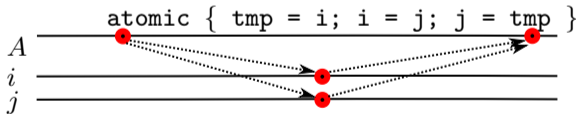
**Program 1**
```
atomic {
  i++;
}
```

**Program 2**
```
atomic {
  j = i;
  i = i+k;
}
```

**Program 3**
```
atomic {
  int tmp = i;
  i = j;
  j = tmp;
}
```

The statements in an `atomic` block execute as *atomic execution*:



- `atomic` only translatable when a corresponding atomic CPU instruction exist
- the notion of requesting *atomic execution* is a general concept

# Wait-Free Synchronization

Wait-Free algorithms are limited to a single instruction:

- no control flow possible, no behavioral change depending on data
- often, there are instructions that execute an operation conditionally

| Program 4 | Program 5 | Program 6 |
|---|---|---|
| ```atomic {    r = b;   b = 0; }``` | ```atomic {    r = b;   b = 1; }``` | ```atomic {    r = (k==i);   if (r) i = j; }``` |

Operations *update* a memory cell and *return* the previous value.

- the first two operations can be seen as setting a flag `b` to $v \in \{0, 1\}$ and returning its previous state.
  - ▶ the operation implementing programs 4 and 5 is called *set-and-test*
- the third case generalizes this to setting a variable `i` to the value of `j`, if `i`'s old value is equal to `k`'s.
  - ▶ the operation implementing program 6 is called *compare-and-swap*

# Wait-Free Synchronization

Wait-Free algorithms are limited to a single instruction:
- no control flow possible, no behavioral change depending on data
- often, there are instructions that execute an operation conditionally

| Program 4 | Program 5 | Program 6 |
|---|---|---|
| ```
atomic {
  r = b;
  b = 0;
}
``` | ```
atomic {
  r = b;
  b = 1;
}
``` | ```
atomic {
  r = (k==i);
  if (r) i = j;
}
``` |

Operations *update* a memory cell and *return* the previous value.
- the first two operations can be seen as setting a flag b to $v \in \{0, 1\}$ and returning its previous state.
  - ▸ the operation implementing programs 4 and 5 is called *set-and-test*
- the third case generalizes this to setting a variable i to the value of j, if i's old value is equal to k's.
  - ▸ the operation implementing program 6 is called *compare-and-swap*
⤳ use as *building blocks* for algorithms that can *fail*

**Lock-Free Algorithms**

If a *wait-free* implementation is not possible, a *lock-free* implementation might still be viable.

# **Lock-Free Algorithms**

If a *wait-free* implementation is not possible, a *lock-free* implementation might still be viable.

Common usage pattern for *compare and swap*:

1. read the initial value in $i$ into $k$ (using memory barriers)
2. compute a new value $j = f(k)$
3. update $i$ to $j$ if $i = k$ still holds
4. go to first step if $i \neq k$ meanwhile

# Lock-Free Algorithms

If a *wait-free* implementation is not possible, a *lock-free* implementation might still be viable.

Common usage pattern for *compare and swap*:

1. read the initial value in $i$ into $k$ (using memory barriers)
2. compute a new value $j = f(k)$
3. update $i$ to $j$ if $i = k$ still holds
4. go to first step if $i \neq k$ meanwhile

⚠ note: $i = k$ must imply that no thread has updated $i$

# Lock-Free Algorithms

If a *wait-free* implementation is not possible, a *lock-free* implementation might still be viable.

Common usage pattern for *compare and swap*:

1. read the initial value in $i$ into $k$ (using memory barriers)
2. compute a new value $j = f(k)$
3. update $i$ to $j$ if $i = k$ still holds
4. go to first step if $i \neq k$ meanwhile

⚠ note: $i = k$ must imply that no thread has updated $i$

## General recipe for *lock-free* algorithms

- given a compare-and-swap operation for $n$ bytes
- try to group variables for which an invariant must hold into $n$ bytes
- read these bytes atomically
- compute a new value
- perform a compare-and-swap operation on these $n$ bytes

# Lock-Free Algorithms

If a *wait-free* implementation is not possible, a *lock-free* implementation might still be viable.

Common usage pattern for *compare and swap*:

1. read the initial value in $i$ into $k$ (using memory barriers)
2. compute a new value $j = f(k)$
3. update $i$ to $j$ if $i = k$ still holds
4. go to first step if $i \neq k$ meanwhile

⚠ note: $i = k$ must imply that no thread has updated $i$

## General recipe for *lock-free* algorithms

- given a compare-and-swap operation for $n$ bytes
- try to group variables for which an invariant must hold into $n$ bytes
- read these bytes atomically
- compute a new value
- perform a compare-and-swap operation on these $n$ bytes

⇝ computing new value must be *repeatable* or *pure*

# **Limitations of Wait- and Lock-Free Algorithms** 

Wait-/Lock-Free algorithms are severely limited in terms of their computation:

- restricted to the semantics of a *single* atomic operation
- set of atomic operations is architecture specific, but often includes
  - ► exchange of a memory cell with a register
  - ► compare-and-swap of a register with a memory cell
  - ► fetch-and-add on integers in memory
  - ► modify-and-test on bits in memory
- provided instructions usually allow only one memory operand

# Limitations of Wait- and Lock-Free Algorithms

Wait-/Lock-Free algorithms are severely limited in terms of their computation:

- restricted to the semantics of a *single* atomic operation
- set of atomic operations is architecture specific, but often includes
  - exchange of a memory cell with a register
  - compare-and-swap of a register with a memory cell
  - fetch-and-add on integers in memory
  - modify-and-test on bits in memory
- provided instructions usually allow only one memory operand

⤳ Lock-Free instructions as *building blocks* for *Locks*

**Locked Atomic Executions**

# Locks

### Definition (Lock)

A lock is a data structure that

- can be *acquired* and *released*
- ensures *mutual exclusion*: only one thread may hold the lock at a time
- *blocks* other threads attempts to acquire while held by a different thread
- protects a *critical section*: a piece of code that may produce incorrect results when entered concurrently from several threads

⚠ may *deadlock* the program

# Semaphores and Mutexes

A (counting) *semaphore* is an integer `s` with the following operations:

```
void signal(int *s) {
  atomic { *s = *s + 1; }
}
```

```
void wait(int *s) {
  bool avail;
  do {
    atomic {
      avail = *s>0;
      if (avail) (*s)--;
    }
  } while (!avail);
}
```

# Semaphores and Mutexes

A (counting) *semaphore* is an integer `s` with the following operations:

```c
void signal(int *s) {
  atomic { *s = *s + 1; }
}
```

```c
void wait(int *s) {
  bool avail;
  do {
    atomic {
      avail = *s>0;
      if (avail) (*s)--;
    }
  } while (!avail);
}
```

A counting semaphore can track how many resources are still available.

- a thread *acquiring* a resource executes `wait()`
- if a resource is still available, `wait()` returns
- once a thread finishes using a resource, it calls `signal()` to *release*

# Semaphores and Mutexes

A (counting) *semaphore* is an integer `s` with the following operations:

```c
void signal(int *s) {
  atomic { *s = *s + 1; }
}
```

```c
void wait(int *s) {
  bool avail;
  do {
    atomic {
      avail = *s>0;
      if (avail) (*s)--;
    }
  } while (!avail);
}
```

A counting semaphore can track how many resources are still available.

- a thread *acquiring* a resource executes `wait()`
- if a resource is still available, `wait()` returns
- once a thread finishes using a resource, it calls `signal()` to *release*

Special case: initializing with $s = 1$ gives a *binary* semaphore:

- can be used to block and unblock a thread
- can be used to protect a single resource

$\rightsquigarrow$ in this case the data structure is also called *mutex*

## Implementation of Semaphores

A *semaphore* does not have to wait busily:

```
void signal(int *s) {
  atomic { *s = *s + 1; }
  wake(s);
}
```

```
void wait(int *s) {
  bool avail;
  do {
    atomic {
      avail = *s>0;
      if (avail) (*s)--;
    }
    if (!avail) de_schedule(s);
  } while (!avail);
}
```

## Implementation of Semaphores

A *semaphore* does not have to wait busily:

```
void signal(int *s) {
  atomic { *s = *s + 1; }
  wake(s);
}
```

```
void wait(int *s) {
  bool avail;
  do {
    atomic {
      avail = *s>0;
      if (avail) (*s)--;
    }
    if (!avail) de_schedule(s);
  } while (!avail);
}
```

Busy waiting is avoided:

- a thread failing to decrease *s executes de_schedule()
- de_schedule() enters the operating system and inserts the current thread into a queue of threads that will be woken up when *s becomes non-zero, usually by *monitoring writes* to s (⤳ *FUTEX_WAIT*)
- once a thread calls wake(s), the first thread $t$ waiting on s is extracted
- the operating system lets $t$ return from its call to de_schedule()

## Practical Implementation of Semaphores

Certain optimisations are possible:

```c
void signal(int *s) {
  atomic { *s = *s + 1; }
  wake(s);
}
```

```c
void wait(int *s) {
  bool avail;
  do {
    atomic {
      avail = *s>0;
      if (avail) (*s)--;
    }
    if (!avail) de_schedule(s);
  } while (!avail);
}
```

In general, the implementation is more complicated
- `wait()` may busy wait for a few iterations
  ▸ avoids de-scheduling if the lock is released frequently
  ▸ better throughput for semaphores that are held for a short time
- `wake(s)` informs the scheduler that `s` has been written to

# Practical Implementation of Semaphores

Certain optimisations are possible:

```c
void signal(int *s) {
  atomic { *s = *s + 1; }
  wake(s);
}
```

```c
void wait(int *s) {
  bool avail;
  do {
    atomic {
      avail = *s>0;
      if (avail) (*s)--;
    }
    if (!avail) de_schedule(s);
  } while (!avail);
}
```

In general, the implementation is more complicated

- `wait()` may busy wait for a few iterations
  ▶ avoids de-scheduling if the lock is released frequently
  ▶ better throughput for semaphores that are held for a short time
- `wake(s)` informs the scheduler that `s` has been written to
↝ using a semaphore with a single core reduces to
  `if (*s) (*s)--; /* critical section */ (*s)++;`

One common use of semaphores is to guarantee mutual exclusion.

$\leadsto$ in this case, a binary semaphore is also called a *mutex*

**e.g.** add a lock to the double-ended queue data structure

⚠ decide what needs protection and what not

# **Monitors: An Automatic, Re-entrant Mutex**

Often, a data structure can be made thread-safe by

- *acquiring* a lock upon *entering* a function of the data structure
- *releasing* the lock upon *exit* from this function

# **Monitors: An Automatic, Re-entrant Mutex**

Often, a data structure can be made thread-safe by

- *acquiring* a lock upon *entering* a function of the data structure
- *releasing* the lock upon *exit* from this function

Locking each procedure body that accesses a data structure:

1. is a re-occurring pattern, should be generalized
2. becomes problematic in recursive calls: it blocks

**E.g.** a thread $t$ waits for a data structure to be filled

- ▸ $t$ will call `pop()` and obtain `-1`
- ▸ $t$ then has to call again, until an element is available
- ⤳ $t$ is busy waiting and produces contention on the lock ⚠

# Monitors: An Automatic, Re-entrant Mutex

Often, a data structure can be made thread-safe by

- *acquiring* a lock upon *entering* a function of the data structure
- *releasing* the lock upon *exit* from this function

Locking each procedure body that accesses a data structure:

1. is a re-occurring pattern, should be generalized
2. becomes problematic in recursive calls: it blocks

**E.g.** a thread $t$ waits for a data structure to be filled

- ▶ $t$ will call `pop()` and obtain `-1`
- ▶ $t$ then has to call again, until an element is available
- ⇝ $t$ is busy waiting and produces contention on the lock ⚠

*Monitor*: a mechanism to address these problems:

1. a procedure associated with a monitor acquires a lock on entry and releases it on exit
2. if that lock is already taken by the current thread, proceed

## Monitors: An Automatic, Re-entrant Mutex

Often, a data structure can be made thread-safe by

- *acquiring* a lock upon *entering* a function of the data structure
- *releasing* the lock upon *exit* from this function

Locking each procedure body that accesses a data structure:

1. is a re-occurring pattern, should be generalized
2. becomes problematic in recursive calls: it blocks

**E.g.** a thread $t$ waits for a data structure to be filled

- ▸ $t$ will call `pop()` and obtain `-1`
- ▸ $t$ then has to call again, until an element is available
- ↝ $t$ is busy waiting and produces contention on the lock ⚠

*Monitor*: a mechanism to address these problems:

1. a procedure associated with a monitor acquires a lock on entry and releases it on exit
2. if that lock is already taken by the current thread, proceed

↝ we need a way to release the lock after the return of the last recursive call

## Implementation of a Basic Monitor

A monitor contains a semaphore `count` and the id `tid` of the occupying thread:

```
typedef struct monitor mon_t;
struct monitor { int tid; int count; };
void monitor_init(mon_t* m) { memset(m, 0, sizeof(mon_t)); }
```

## Implementation of a Basic Monitor

A monitor contains a semaphore count and the id tid of the occupying thread:

```
typedef struct monitor mon_t;
struct monitor { int tid; int count; };
void monitor_init(mon_t* m) { memset(m, 0, sizeof(mon_t)); }
```

Define monitor_enter and monitor_leave:

- ensure mutual exclusion of accesses to mon_t
- track how many times we called a monitored procedure recursively

```
void monitor_enter(mon_t *m) {
 bool mine = false;
 while (!mine) {
   mine = thread_id()==m->tid;
   if (mine) m->count++; else
   atomic {
     if (m->tid==0) {
       m->tid = thread_id();
       mine = true; m->count=1;
   } };
   if (!mine) de_schedule(&m->tid);
} }
```

```
void monitor_leave(mon_t *m) {
  m->count--;
  if (m->count==0) {
    atomic {
      m->tid=0;
    }
    wake(&m->tid);
  }
}
```

# Condition Variables

✓ Monitors simplify the construction of thread-safe resources.

Still: Efficiency problem when using resource to synchronize:

**E.g.** a thread $t$ waits for a data structure to be filled:

- ▶ $t$ will call `pop()` and obtain `-1`
- ▶ $t$ then has to call again, until an element is available
- ⤳ $t$ is busy waiting and produces contention on the lock

# Condition Variables

$\checkmark$ Monitors simplify the construction of thread-safe resources.

Still: Efficiency problem when using resource to synchronize:

**E.g.** a thread $t$ waits for a data structure to be filled:

- ▸ $t$ will call `pop()` and obtain `-1`
- ▸ $t$ then has to call again, until an element is available
- ⤳ $t$ is busy waiting and produces contention on the lock

Idea: create a *condition variable* on which to block while waiting:

```
struct monitor { int tid; int count; int cond; int cond2;... };
```

# Condition Variables

✓ Monitors simplify the construction of thread-safe resources.

Still: Efficiency problem when using resource to synchronize:

**E.g.** a thread $t$ waits for a data structure to be filled:

- $t$ will call `pop()` and obtain `-1`
- $t$ then has to call again, until an element is available
- ⤳ $t$ is busy waiting and produces contention on the lock

Idea: create a *condition variable* on which to block while waiting:

```
struct monitor { int tid; int count; int cond; int cond2;... };
```

Define these two functions:

1. `wait` for the condition to become true
   - called while being *inside* the monitor
   - temporarily *releases* the monitor and blocks
   - when *signalled*, re-acquires the monitor and returns

2. `signal` waiting threads that they may be able to proceed
   - one/all waiting threads that called *wait* will be woken up, two possibilities:
     **signal-and-urgent-wait** : the *signalling* thread suspends and continues once the *signalled* thread has released the monitor
     **signal-and-continue** the *signalling* thread continues, any *signalled* thread enters when the monitor becomes available

# Signal-And-Urgent-Wait Semantics

Requires one queue for each condition $c$ and a suspended queue $s$:



- a thread who tries to enter a monitor is added to queue $e$ if the monitor is occupied
- a call to `wait` on condition $a$ adds thread to the queue $a.q$
- a call to `signal` for $a$ adds thread to queue $s$ (suspended)
- one thread from the $a$ queue is woken up
- `signal` on $a$ is a no-op if $a.q$ is empty
- if a thread leaves, it wakes up one thread waiting on $s$
- if $s$ is empty, it wakes up one thread from $e$

# Signal-And-Urgent-Wait Semantics

Requires one queue for each condition $c$ and a suspended queue $s$:



source: http://en.wikipedia.org/wiki/Monitor_(synchronization)

- a thread who tries to enter a monitor is added to queue $e$ if the monitor is occupied
- a call to `wait` on condition $a$ adds thread to the queue $a.q$
- a call to `signal` for $a$ adds thread to queue $s$ (suspended)
- one thread from the $a$ queue is woken up
- `signal` on $a$ is a no-op if $a.q$ is empty
- if a thread leaves, it wakes up one thread waiting on $s$
- if $s$ is empty, it wakes up one thread from $e$

$\rightsquigarrow$ queue $s$ has priority over $e$

# Signal-And-Continue Semantics

Here, the signal function is usually called notify.



- a call to wait on condition $a$ adds thread to the queue $a.q$
- a call to notify for $a$ adds one thread from $a.q$ to $e$ (unless $a.q$ is empty)
- if a thread leaves, it wakes up one thread waiting on $e$

# Signal-And-Continue Semantics

Here, the `signal` function is usually called `notify`.



- a call to `wait` on condition $a$ adds thread to the queue $a.q$
- a call to `notify` for $a$ adds one thread from $a.q$ to $e$ (unless $a.q$ is empty)
- if a thread leaves, it wakes up one thread waiting on $e$

$\rightsquigarrow$ signalled threads compete for the monitor

- assuming FIFO ordering on $e$, threads who tried to enter between `wait` and `notify` will run first
- need additional queue $s$ if waiting threads should have priority

## Implementing Condition Variables

We implement the simpler *signal-and-continue* semantics for a single condition variable:

⇝ a *notified* thread is simply woken up and competes for the monitor

```
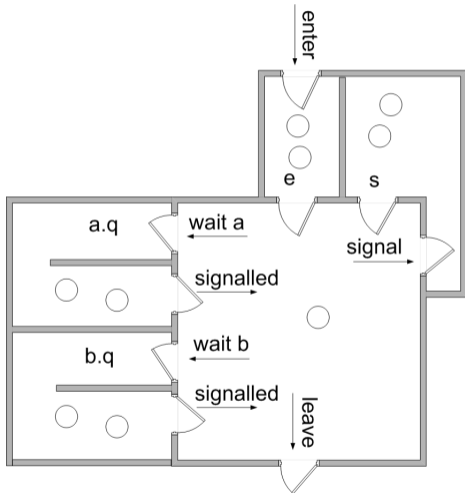void cond_wait(mon_t *m) {
  assert(m->tid==thread_id());
  int old_count = m->count;
  m->tid = 0;
  wait(&m->cond);
  bool next_to_enter;
  do {
    atomic {
      next_to_enter = m->tid==0;
      if (next_to_enter) {
        m->tid = thread_id();
        m->count = old_count;
      }
    }
    if (!next_to_enter) de_schedule(&m->tid);
  } while (!next_to_enter);}
```

```
void cond_notify(mon_t *m) {
  // wake up other threads
  signal(&m->cond);
}
```

# A Note on Notify

With *signal-and-continue* semantics, two notify functions exist:

1. `notify`: wakes up exactly one thread waiting on condition variable
2. `notifyAll`: wakes up all threads waiting on a condition variable

# A Note on Notify

With *signal-and-continue* semantics, two notify functions exist:

1. `notify`: wakes up exactly one thread waiting on condition variable
2. `notifyAll`: wakes up all threads waiting on a condition variable

⚠ an implementation often becomes easier if `notify` means *notify some*

⤳ programmer should assume that thread is not the only one woken up

## Monitors with a Single Condition Variable

Monitors with a single condition variable are built into *Java* and *C#*:



source: http://en.wikipedia.org/wiki/Monitor_(synchronization)

```
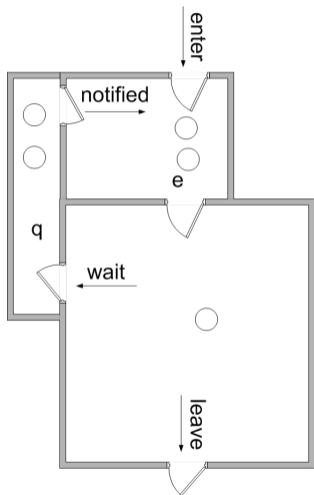class C {
  public synchronized void f() {
    // body of f
}}
```

is equivalent to

```
class C {
  public void f() {
    monitor_enter(this);
    // body of f
    monitor_leave(this);
}}
```

with `Object` containing:

```
private int mon_var;
private int mon_count;
private int cond_var;
protected void monitor_enter();
protected void monitor_leave();
```

**Deadlocks**

**Definition (Deadlock)**

A deadlock is a situation in which two processes are waiting for the respective other to finish, and thus neither ever does.

(The definition generalizes to a set of actions with a cyclic dependency.)

### Definition (Deadlock)

A deadlock is a situation in which two processes are waiting for the respective other to finish, and thus neither ever does.

(The definition generalizes to a set of actions with a cyclic dependency.)

Consider this Java class:

```java
class Foo {
  public Foo other = null;
  public synchronized void bar() {
    ... if (*) other.bar(); ...
  }
}
```

and two instances:

```java
Foo a = new Foo(), b = new Foo();
a.other = b; b.other = a;
// in parallel:
a.bar() || b.bar();
```

Sequence leading to a deadlock:

- threads $A$ and $B$ execute `a.bar()` and `b.bar()`
- `a.bar()` acquires the monitor of `a`
- `b.bar()` acquires the monitor of `b`
- $A$ happens to execute `other.bar()`
- $A$ blocks on the monitor of $b$
- $B$ happens to execute `other.bar()`
- ⤳ both *block* indefinitely

## **Deadlocks with Monitors**

### **Definition (Deadlock)**

A deadlock is a situation in which two processes are waiting for the respective other to finish, and thus neither ever does.

(The definition generalizes to a set of actions with a cyclic dependency.)

Consider this Java class:

```java
class Foo {
  public Foo other = null;
  public synchronized void bar() {
    ... if (*) other.bar(); ...
  }
}
```

and two instances:

```java
Foo a = new Foo(), b = new Foo();
a.other = b; b.other = a;
// in parallel:
a.bar() || b.bar();
```

Sequence leading to a deadlock:

- threads $A$ and $B$ execute `a.bar()` and `b.bar()`
- `a.bar()` acquires the monitor of `a`
- `b.bar()` acquires the monitor of `b`
- $A$ happens to execute `other.bar()`
- $A$ blocks on the monitor of $b$
- $B$ happens to execute `other.bar()`
- ⤳ both *block* indefinitely

How can this situation be avoided?

# Treatment of Deadlocks

Observation: Deadlocks occur if the following four conditions hold
[Coffman et al.(1971)Coffman, Elphick, and Shoshani]:

1. *mutual exclusion*: processes require exclusive access
2. *wait for*: a process holds resources while waiting for more
3. *no preemption*: resources cannot be taken away form processes
4. *circular wait*: waiting processes form a cycle

## Treatment of Deadlocks

Observation: Deadlocks occur if the following four conditions hold
[Coffman et al.(1971)Coffman, Elphick, and Shoshani]:

1. *mutual exclusion*: processes require exclusive access
2. *wait for*: a process holds resources while waiting for more
3. *no preemption*: resources cannot be taken away form processes
4. *circular wait*: waiting processes form a cycle

The occurrence of deadlocks can be:

1. *ignored*: for the lack of better approaches, can be reasonable if deadlocks are rare
2. *detection*: check within OS for a cycle, requires ability to *preempt*
3. prevention: design programs to be deadlock-free
4. *avoidance*: use additional information about a program that allows the OS to schedule threads so that they do not deadlock

# **Treatment of Deadlocks**

Observation: Deadlocks occur if the following four conditions hold
[Coffman et al.(1971)Coffman, Elphick, and Shoshani]:

1. *mutual exclusion*: processes require exclusive access
2. *wait for*: a process holds resources while waiting for more
3. *no preemption*: resources cannot be taken away form processes
4. *circular wait*: waiting processes form a cycle

The occurrence of deadlocks can be:

1. *ignored*: for the lack of better approaches, can be reasonable if deadlocks are rare
2. *detection*: check within OS for a cycle, requires ability to *preempt*
3. prevention: design programs to be deadlock-free
4. *avoidance*: use additional information about a program that allows the OS to schedule threads so that they do not deadlock

$\leadsto$ *prevention* is the only safe approach on standard operating systems

- can be achieved using *lock-free* algorithms
- but what about algorithms that require locking?

# Deadlock Prevention through Partial Order

Observation: A cycle cannot occur if locks are *partially ordered*.

**Definition (lock sets)**

Let $L$ denote the set of locks. We call $\lambda(p) \subseteq L$ the lock set at $p$, i.e. the set of locks that may be in the "acquired" state at program point $p$.

# Deadlock Prevention through Partial Order

Observation: A cycle cannot occur if locks are *partially ordered*.

**Definition (lock sets)**

Let $L$ denote the set of locks. We call $\lambda(p) \subseteq L$ the lock set at $p$, i.e. the set of locks that may be in the "acquired" state at program point $p$.

We require the transitive closure $\sigma^+$ of a relation $\sigma$:

**Definition (transitive closure)**

Let $\sigma \subseteq X \times X$ be a relation. Its transitive closure is $\sigma^+ = \bigcup_{i \in \mathbb{N}} \sigma^i$ where

$$
\begin{aligned}
\sigma^0 &= \sigma \\
\sigma^{i+1} &= \{\langle x_1, x_3 \rangle \mid \exists x_2 \in X . \langle x_1, x_2 \rangle \in \sigma^i \wedge \langle x_2, x_3 \rangle \in \sigma^i\} \cup \sigma^i
\end{aligned}
$$

# Deadlock Prevention through Partial Order

Observation: A cycle cannot occur if locks are *partially ordered*.

**Definition (lock sets)**

Let $L$ denote the set of locks. We call $\lambda(p) \subseteq L$ the lock set at $p$, i.e. the set of locks that may be in the "acquired" state at program point $p$.

We require the transitive closure $\sigma^+$ of a relation $\sigma$:

**Definition (transitive closure)**

Let $\sigma \subseteq X \times X$ be a relation. Its transitive closure is $\sigma^+ = \bigcup_{i \in \mathbb{N}} \sigma^i$ where

$$\sigma^0 = \sigma$$
$$\sigma^{i+1} = \{\langle x_1, x_3 \rangle \mid \exists x_2 \in X . \langle x_1, x_2 \rangle \in \sigma^i \wedge \langle x_2, x_3 \rangle \in \sigma^i\} \cup \sigma^i$$

Each time a lock is acquired, we track the lock set at $p$:

**Definition (lock order)**

Define $\lhd \subseteq L \times L$ such that $l \lhd l'$ iff $l \in \lambda(p)$ and the statement at $p$ is of the form `wait(l')` or `monitor_enter(l')`. Define the lock order $\prec \; = \; \lhd^+$.

# **Freedom of Deadlock**

The following holds for a program with mutexes and monitors:

**Theorem (freedom of deadlock)**

*If there exists no $a \in L$ with $a \prec a$ then the program is free of deadlocks.*

## Freedom of Deadlock

The following holds for a program with mutexes and monitors:

**Theorem (freedom of deadlock)**

*If there exists no $a \in L$ with $a \prec a$ then the program is free of deadlocks.*

Suppose a program blocks on semaphores (mutexes) $L_S$ and on monitors $L_M$ such that $L = L_S \cup L_M$.

**Theorem (freedom of deadlock for monitors)**

*If $\forall a \in L_S \,.\, a \not\prec a$ and $\forall a \in L_M, b \in L \,.\, a \prec b \wedge b \prec a \Rightarrow a = b$ then the program is free of deadlocks.*

## **Freedom of Deadlock**

The following holds for a program with mutexes and monitors:

**Theorem (freedom of deadlock)**

*If there exists no $a \in L$ with $a \prec a$ then the program is free of deadlocks.*

Suppose a program blocks on semaphores (mutexes) $L_S$ and on monitors $L_M$ such that $L = L_S \cup L_M$.

**Theorem (freedom of deadlock for monitors)**

*If $\forall a \in L_S \, . \, a \nprec a$ and $\forall a \in L_M, b \in L \, . \, a \prec b \wedge b \prec a \Rightarrow a = b$ then the program is free of deadlocks.*

Note: the set $L$ contains *instances* of a lock.

- the set of lock instances can vary at runtime
- if we statically want to ensure that deadlocks cannot occur:
  - ▶ summarize every lock/monitor that may have several instances into one
  - ▶ a summary lock/monitor $\bar{a} \in L_M$ represents several concrete ones
  - ▶ thus, if $\bar{a} \prec \bar{a}$ then this might not be a self-cycle
- ⤳ require that $\bar{a} \nprec \bar{a}$ for all summarized monitors $\bar{a} \in L_M$

⚠ fix a representation for locksets

↝ in our case: $L$ comprises all lines, where any object is created.

```
0:  Foo  a = new Foo();
1:  Foo  b = new Foo();
2:  a.other = b;
3:  b.other = a;
4:
5:
6:  bar(&a); || bar(&b);
7:
```

```
8:  void bar(this) {
9:    monitor_enter(this);
10:   if (*) {
11:      ...
12:      bar(&other);
13:      ...
14:   }
15:   monitor_leave(this);
16: }
```

| Lockorder | ◁ | |

# Inferring locksets and lockset order in practice

⚠ fix a representation for locksets

⤳ in our case: $L$ comprises all lines, where any object is created.

```
0:  Foo  a = new Foo();
1:  Foo  b = new Foo();
2:  a.other = b;
3:  b.other = a;
4:
5:
6:  bar(&a); || bar(&b);
7:
```

$\lambda(8) = \{\}$

```
8:   void bar(this) {
9:     monitor_enter(this);
10:    if (*) {
11:       ...
12:       bar(&other);
13:       ...
14:    }
15:    monitor_leave(this);
16:  }
```

| Lockorder | ◁ | |
|-----------|---|--|

# Inferring locksets and lockset order in practice

⚠ fix a representation for locksets

⤳ in our case: $L$ comprises all lines, where any object is created.

```
0:  Foo  a = new Foo();
1:  Foo  b = new Foo();
2:  a.other = b;
3:  b.other = a;
4:
5:
6:  bar(&a); || bar(&b);
7:
```

$\lambda(9) = \{l_0, l_1\}$

```
8:  void bar(this) {
9:    monitor_enter(this);
10:   if (*) {
11:     ...
12:     bar(&other);
13:     ...
14:   }
15:   monitor_leave(this);
16:  }
```

$this = \{\&a, \&b\}$

| Lockorder | ◁ | |
|---|---|---|

# Inferring locksets and lockset order in practice

⚠ fix a representation for locksets

⤳ in our case: $L$ comprises all lines, where any object is created.

```
0:  Foo  a = new Foo();
1:  Foo  b = new Foo();
2:  a.other = b;
3:  b.other = a;
4:
5:
6:  bar(&a); || bar(&b);
7:
```

```
8:   void bar(this) {
9:     monitor_enter(this);
10:    if (*) {
11:      ...
12:      bar(&other);
13:      ...
14:    }
15:    monitor_leave(this);
16:  }
```

$\texttt{this} = \{\&a, \&b\}$

$\lambda(11) = \{l_0, l_1\}$

| Lockorder | ◁ | |
|-----------|---|---|

# Inferring locksets and lockset order in practice

⚠ fix a representation for locksets

⤳ in our case: $L$ comprises all lines, where any object is created.

```
0:  Foo  a = new Foo();
1:  Foo  b = new Foo();
2:  a.other = b;
3:  b.other = a;
4:
5:
6:  bar(&a); || bar(&b);
7:
```

```
8:   void bar(this) {
9:     monitor_enter(this);
10:    if (*) {
11:      ...
12:      bar(&other);
13:      ...
14:    }
15:    mon...
16:  }
```

this $= \{\&a, \&b\}$

$\lambda(11) = \{l_0, l_1\}$

other $= \{\&a, \&b\}$

| Lockorder | ◁ | |
|---|---|---|

⚠ fix a representation for locksets

⤳ in our case: $L$ comprises all lines, where any object is created.

this $= \{\&a, \&b\}$

```
0: Foo  a = new Foo();
1: Foo  b = new Foo();
2: a.other = b;
3: b.other = a;
4:
5:
6: bar(&a); || bar(&b);
7:
```

$\lambda(8) = \{l_0, l_1\}$

```
8:  void bar(this) {
9:    monitor_enter(this);
10:   if (*) {
11:     ...
12:     bar(&other);
13:     ...
14:   }
15:   mon...
16: }
```

other $= \{\&a, \&b\}$

| Lockorder | ◁ | |
|-----------|---|---|

# Inferring locksets and lockset order in practice

⚠ fix a representation for locksets

↝ in our case: $L$ comprises all lines, where any object is created.

```
0:  Foo  a = new Foo();
1:  Foo  b = new Foo();
2:  a.other = b;
3:  b.other = a;
4:
5:
6:  bar(&a);  ||  bar(&b);
7:
```

$\lambda(9) = \{\mathtt{l}_0, \mathtt{l}_1\}$

```
8:  void bar(this) {
9:    monitor_enter(this);
10:   if (*) {
11:     ...
12:     bar(&other);
13:     ...
14:   }
15:   mon...
16: }
```

$\mathtt{this} = \{\&a, \&b\}$

$\mathtt{other} = \{\&a, \&b\}$

| Lockorder | ◁ | $\langle l_0, l_1 \rangle, \langle l_1, l_0 \rangle$ |

# Avoiding Deadlocks in Practice

⚠ What to do when the lock order contains a cycle?

- determining which locks may be acquired at each program point is undecidable
  ↝ lock sets are an approximation
- an array of locks in $L_S$: lock in increasing array index sequence
- if $l \in \lambda(P)$ exists $l' \prec l$ is to be acquired
  ↝ change program: release $l$, acquire $l'$, then acquire $l$ again
  ⚠ inefficient
- if a lock set contains a summarized lock $\bar{a}$ and $\bar{a}$ is to be acquired, we're stuck

**Locks Roundup**

## Atomic Execution and Locks

Consider replacing the specific locks with `atomic` annotations:

**stack: removal**

```
void pop() {
  ...
  wait(&q->t);
  ...
  if (*) { signal(&q->t); return; }
  ...
  if (c) wait(&q->s);
  ...
  if (c) signal(&q->s);
  signal(&q->t);
}
```

## Atomic Execution and Locks

Consider replacing the specific locks with `atomic` annotations:

**stack: removal**

```
void pop() {
  ...
  wait(&q->t);
  ...
  if (*) { signal(&q->t); return; }
  ...
  if (c) wait(&q->s);
  ...
  if (c) signal(&q->s);
  signal(&q->t);
}
```

- nested `atomic` blocks still describe one atomic execution
- $\leadsto$ locks convey additional information over `atomic`
- locks cannot easily be recovered from `atomic` declarations

Writing `atomic` annotations around sequences of statements is a convenient way of programming.

Writing `atomic` annotations around sequences of statements is a convenient way of programming.

*Idea of mutexes:* Implement `atomic` sections with locks:

- a single lock could be used to protect all `atomic` blocks
- more concurrency is possible by using several locks
- some statements might modify variables that are never read by other threads ⇝ no lock required
- statements in one `atomic` block might access variables in a different order to another `atomic` block ⇝ deadlock possible with locks implementation
- creating too many locks can decrease the performance, especially when required to release locks in $\lambda(l)$ when acquiring $l$

## **Outlook**

Writing `atomic` annotations around sequences of statements is a convenient way of programming.

*Idea of mutexes:* Implement `atomic` sections with locks:

- a single lock could be used to protect all `atomic` blocks
- more concurrency is possible by using several locks
- some statements might modify variables that are never read by other threads ⤳ no lock required
- statements in one `atomic` block might access variables in a different order to another `atomic` block ⤳ deadlock possible with locks implementation
- creating too many locks can decrease the performance, especially when required to release locks in $\lambda(l)$ when acquiring $l$

⤳ creating locks automatically is non-trivial and, thus, not standard in programming languages

# Concurrency across Languages

In most systems programming languages (C,C++) we have

- the ability to use *atomic* operations
- ⤳ we can implement *wait-free* algorithms

# Concurrency across Languages

In most systems programming languages (C,C++) we have

- the ability to use *atomic* operations
- $\rightsquigarrow$ we can implement *wait-free* algorithms

In Java, C# and other higher-level languages

- provide monitors and possibly other concepts
- often simplify the programming but incur the same problems

# Concurrency across Languages

In most systems programming languages (C,C++) we have

- the ability to use *atomic* operations
$\leadsto$ we can implement *wait-free* algorithms

In Java, C# and other higher-level languages

- provide monitors and possibly other concepts
- often simplify the programming but incur the same problems

| language | barriers | wait-/lock-free | semaphore | mutex | monitor |
|----------|----------|-----------------|-----------|-------|---------|
| C,C++    | ✓        | ✓               | ✓         | ✓     | (a)     |
| Java,C#  | -        | (b)             | (c)       | ✓     | ✓       |

**(a)** some pthread implementations allow a *reentrant* attribute

**(b)** newer API extensions ( `java.util.concurrent.atomic.*` and `System.Threading.Interlocked` resp.)

**(c)** simulate semaphores using an object with two `synchronized` methods

## Summary

Classification of concurrency algorithms:

- wait-free, lock-free, locked
- next on the agenda: transactional

*Wait-free* algorithms:

- never block, always succeed, never deadlock, no starvation
- very limited in expressivity

*Lock-free* algorithms:

- never block, may fail, never deadlock, may starve
- invariant may only span a few bytes (8 on Intel)

*Locking* algorithms:

- can guard arbitrary code
- can use several locks to enable more fine grained concurrency
- may deadlock
- semaphores are not re-entrant, monitors are

⇝ use algorithm that is best fit

# References

📕 E. G. Coffman, M. Elphick, and A. Shoshani.
System deadlocks.
*ACM Comput. Surv.*, 3(2):67–78, June 1971.
ISSN 0360-0300.

📕 T. Harris, J. Larus, and R. Rajwar.
Transactional memory, 2nd edition.
*Synthesis Lectures on Computer Architecture*, 5(1):1–263, 2010.

# **Programming Languages**

Concurrency: Transactions

Dr. Michael Petter
Winter term 2019

# Abstraction and Concurrency

Two fundamental concepts to build larger software are:

**abstraction** : an object storing certain data and providing certain functionality may be used without reference to its internals

**composition** : several objects can be combined to a new object without interference

Both, *abstraction* and *composition* are closely related, since the ability to compose depends on the ability to abstract from details.

## Abstraction and Concurrency

Two fundamental concepts to build larger software are:

**abstraction** : an object storing certain data and providing certain functionality may be used without reference to its internals

**composition** : several objects can be combined to a new object without interference

Both, *abstraction* and *composition* are closely related, since the ability to compose depends on the ability to abstract from details.

Consider an example:

- a linked list data structure exposes a fixed set of operations to modify the list structure, such as `push()` and `forAll()`
- a set object may internally use the list object and expose a set of operations, including `push()`

The `insert()` operations uses the `forAll()` operation to check if the element already exists and uses `push()` if not.

## Abstraction and Concurrency

Two fundamental concepts to build larger software are:

**abstraction** : an object storing certain data and providing certain functionality may be used without reference to its internals

**composition** : several objects can be combined to a new object without interference

Both, *abstraction* and *composition* are closely related, since the ability to compose depends on the ability to abstract from details.

Consider an example:

- a linked list data structure exposes a fixed set of operations to modify the list structure, such as `push()` and `forAll()`
- a set object may internally use the list object and expose a set of operations, including `push()`

The `insert()` operations uses the `forAll()` operation to check if the element already exists and uses `push()` if not.

Wrapping the linked list in a mutex does not help to make the *set* thread-safe.

⤳ wrap the two calls in `insert()` in a mutex

- but other list operations can still be called ⤳ use the *same* mutex

## Abstraction and Concurrency

Two fundamental concepts to build larger software are:

**abstraction** : an object storing certain data and providing certain functionality may be
used without reference to its internals

**composition** : several objects can be combined to a new object without interference

Both, *abstraction* and *composition* are closely related, since the ability to compose
depends on the ability to abstract from details.

Consider an example:

- a linked list data structure exposes a fixed set of operations to modify the list structure,
  such as `push()` and `forAll()`
- a set object may internally use the list object and expose a set of operations, including
  `push()`

The `insert()` operations uses the `forAll()` operation to check if the element already
exists and uses `push()` if not.

Wrapping the linked list in a mutex does not help to make the *set* thread-safe.

⤳ wrap the two calls in `insert()` in a mutex

- but other list operations can still be called ⤳ use the *same* mutex

⤳ unlike sequential algorithms, thread-safe algorithms cannot always be composed to
give new thread-safe algorithms

Idea: automatically convert `atomic` blocks into code that ensures atomic execution of the statements.

```
atomic {
  // code
  if (cond) retry;
  atomic {
    // more code
  }
  // code
}
```

# Transactional Memory [2]

Idea: automatically convert `atomic` blocks into code that ensures atomic execution of the statements.

```
atomic {
  // code
  if (cond) retry;
  atomic {
    // more code
  }
  // code
}
```

Execute code as *transaction*:

- execute the code of an atomic block
- nested atomic blocks act like a single atomic block
- check that it runs without *conflicts* due to accesses from another thread
- if another thread interferes through conflicting updates:
  - ▸ undo the computation done so far
  - ▸ re-start the transaction
- provide a `retry` keyword similar to the `wait` of monitors

# Semantics of Transactions

The goal is to use transactions to specify *atomic executions*.

Transactions are rooted in databases where they have the *ACID* properties:

# Semantics of Transactions

The goal is to use transactions to specify *atomic executions*.

Transactions are rooted in databases where they have the *ACID* properties:

**atomicity** : a transaction completes or seems not to have run

$\rightsquigarrow$ we call this *failure atomicity* to distinguish it from *atomic executions*

**consistency** : each transaction transforms a consistent state to another consistent state

- a consistent state is one in which certain *invariants* hold
- invariants depend on the application

**isolation** : among each other, transactions do not interfere

$\rightsquigarrow$ coexisting with non-transactional memory, isolation is not so evident

**durability** : the effects are permanent (w.r.t. main memory ✓ )

# **Semantics of Transactions**

The goal is to use transactions to specify *atomic executions*.

Transactions are rooted in databases where they have the *ACID* properties:

**atomicity** : a transaction completes or seems not to have run

$\rightsquigarrow$ we call this *failure atomicity* to distinguish it from *atomic executions*

**consistency** : each transaction transforms a consistent state to another consistent state

- a consistent state is one in which certain *invariants* hold
- invariants depend on the application

**isolation** : among each other, transactions do not interfere

$\rightsquigarrow$ coexisting with non-transactional memory, isolation is not so evident

**durability** : the effects are permanent (w.r.t. main memory $\checkmark$ )

---

**Definition (Semantics of Transactions)**

The result of running concurrent transactions must be identical to *one* execution of them in sequence. ($\rightsquigarrow$ Serialization)

# Consistency During Transactions

**Consistency during a transaction.**

ACID states how committed transactions behave but not what may happen until a transaction commits.

- a transaction, run on an inconsistent state may continue yielding inconsistent states
  ↝ zombie transaction
- in the best case, the zombie transaction will be aborted eventually
- but transactions may cause havoc when run on inconsistent states

```
atomic {                                    // preserved invariant: x==y
  int tmp1 = x;                             atomic {
  int tmp2 = y;                               x = 10;
  assert(tmp1-tmp2==0);                       y = 10;
}                                            }
```

⚠ critical for null pointer derefs or divisions by zero, e.g.

**Definition (opacity)**

A TM system provides *opacity* if failing transactions are serializable w.r.t. committing transactions.

↝ failing transactions still see a consistent view of memory

# **Weak- and Strong Isolation**

Can we mix transactions with code accessing memory non-transactionally?

- *strong isolation* retains order between accesses to TM and non-TM
- In *weak isolation*, guarantees are only given about memory accessed inside `atomic`

# Weak- and Strong Isolation

Can we mix transactions with code accessing memory non-transactionally?

- *strong isolation* retains order between accesses to TM and non-TM
- In *weak isolation*, guarantees are only given about memory accessed inside `atomic`
  - ▸ no conflict detection for non-transactional accesses
  - ▸ ⚠ standard *race problems*, e.g.

```
// Thread 1
atomic {
  x = 42;
}
```

```
// Thread 2
int tmp = x;
```

⤳ give programs with races the same semantics as if using a single global lock for all `atomic` blocks

---

### Definition (SLA)

The *single-lock atomicity* is a model in which the program executes as if all transactions acquire a single, program-wide mutual exclusion lock.

# Weak- and Strong Isolation

Can we mix transactions with code accessing memory non-transactionally?

- *strong isolation* retains order between accesses to TM and non-TM
- In *weak isolation*, guarantees are only given about memory accessed inside `atomic`
  - ▸ no conflict detection for non-transactional accesses
  - ▸ ⚠ standard *race problems*, e.g.

```
// Thread 1
atomic {
  x = 42;
}
```
```
// Thread 2
int tmp = x;
```

↝ give programs with races the same semantics as if using a single global lock for all `atomic` blocks

---

### Definition (SLA)

The *single-lock atomicity* is a model in which the program executes as if all transactions acquire a single, program-wide mutual exclusion lock.

---

↝ like *sequential consistency*, SLA is a statement about program equivalence

## Disadvantages of the SLA model

The SLA model is *simple* but often too strong:

**1** SLA has a weaker *progress* guarantee than a transaction should have

```
// Thread 1
atomic {
  while (true) {};
}
```

```
// Thread 2
atomic {
  int tmp = x; // x in TM
}
```

**2** SLA correctness is too strong in practice

```
// Thread 1
data = 1;
atomic {
}
ready = 1;
```

```
// Thread 2
atomic {
  int tmp = data;
  // Thread 1 not in atomic
  if (ready) {
    // use tmp
  }
}
```

  ▸ under the SLA model, `atomic {}` acts as barrier
  ▸ intuitively, the two transactions should be independent rather than synchronize

⤳ need a weaker model for more flexible implementation of *strong isolation*

# Transactional Sequential Consistency

How about a more permissive view of transaction semantics?

- TM should not have the blocking behaviour of locks
- $\rightsquigarrow$ the programmer cannot rely on synchronization

## Definition (TSC)

The *transactional sequential consistency* is a model in which the accesses within each transaction are sequentially consistent.

# Transactional Sequential Consistency

How about a more permissive view of transaction semantics?

- TM should not have the blocking behaviour of locks
- ⤳ the programmer cannot rely on synchronization

---

**Definition (TSC)**

The *transactional sequential consistency* is a model in which the accesses within each transaction are sequentially consistent.



```
                    atomic { k = i+j; }
A
i
j
k
B
     atomic { k = i+j; }              k=42
```

- TSC is weaker: gives *strong isolation*, but allows parallel execution ✓
- TSC is stronger: accesses within a transaction may *not* be re-ordered ⚠

# Transactional Sequential Consistency

How about a more permissive view of transaction semantics?

- TM should not have the blocking behaviour of locks
$\rightsquigarrow$ the programmer cannot rely on synchronization

---

**Definition (TSC)**

The *transactional sequential consistency* is a model in which the accesses within each transaction are sequentially consistent.



- TSC is weaker: gives *strong isolation*, but allows parallel execution ✓
- TSC is stronger: accesses within a transaction may *not* be re-ordered ⚠
$\rightsquigarrow$ actual implementations use TSC with some *race free* re-orderings

**Software Transactional Memory**

# Translation of `atomic`-Blocks

A TM system must track which shared memory locations are accessed:
- convert every read access `x` from a shared variable to `ReadTx(&x)`
- convert every write access `x=e` to a shared variable to `WriteTx(&x,e)`

Convert `atomic` blocks as follows:

```
atomic {
  // code
}
```

$\Longrightarrow$

```
do {
  StartTx();
  // code with ReadTx and WriteTx
} while (!CommitTx());
```

# **Translation of `atomic`-Blocks**

A TM system must track which shared memory locations are accessed:
- convert every read access `x` from a shared variable to `ReadTx(&x)`
- convert every write access `x=e` to a shared variable to `WriteTx(&x,e)`

Convert `atomic` blocks as follows:

```
atomic {
  // code
}
```
$\implies$
```
do {
  StartTx();
  // code with ReadTx and WriteTx
} while (!CommitTx());
```

- translation can be done using a pre-processor
  - determining a minimal set of memory accesses that need to be transactional requires a good static analysis
  - *idea*: translate all accesses to global variables and the heap as TM
  - more fine-grained control using manual translation
- an actual implementation might provide a `retry` keyword
  - when executing `retry`, the transaction aborts and re-starts
  - the transaction will again wind up at `retry` unless its *read set* changes
- $\rightsquigarrow$ block until a variable in the read-set has changed
  - similar to condition variables in monitors $\checkmark$

# A Software TM Implementation

A software TM implementation allocates a *transaction descriptor* to store data specific to each `atomic` block, for instance:

- *undo-log* of all writes which have to be undone if a commit fails
- *redo-log* of all writes which are postponed until a commit
- *read-* and *write-set*: locations accessed so far
- *read-* and *write-version*: time stamp when value was accessed

# A Software TM Implementation

A software TM implementation allocates a *transaction descriptor* to store data specific to each `atomic` block, for instance:

- *undo-log* of all writes which have to be undone if a commit fails
- *redo-log* of all writes which are postponed until a commit
- *read-* and *write-set*: locations accessed so far
- *read-* and *write-version*: time stamp when value was accessed

## Example:

Consider the TL2 STM (software transactional memory) implementation [1]:

- provides *opacity*: zombie transactions do not see inconsistent state
- uses *lazy versioning*: writes are stored in a *redo*-log and done on commit
- *validating conflict detection*: accessing a modified address aborts

The idea: obtain a version from the global counter on starting the transaction, the *read-version*, and watch out for accesses to newer versions throughout the transaction.

# Principles of TL2

TUM

The idea: obtain a version from the global counter on starting the transaction, the *read-version*, and watch out for accesses to newer versions throughout the transaction.

- A read `ReadTx` from a field at `offset` of object `obj` aborts,
  - when the objects version is younger than the transaction
  - when the object is locked at the moment of access

  or returns the read value and adds the accessed memory address to the *read-set*.

The idea: obtain a version from the global counter on starting the transaction, the *read-version*, and watch out for accesses to newer versions throughout the transaction.

- A read `ReadTx` from a field at `offset` of object `obj` aborts,
  - ▶ when the objects version is younger than the transaction
  - ▶ when the object is locked at the moment of access

  or returns the read value and adds the accessed memory address to the *read-set*.
- `WriteTx` is simpler: add or update the location in the *redo-log*.

The idea: obtain a version from the global counter on starting the transaction, the *read-version*, and watch out for accesses to newer versions throughout the transaction.

- A read `ReadTx` from a field at `offset` of object `obj` aborts,
  - when the objects version is younger than the transaction
  - when the object is locked at the moment of access

  or returns the read value and adds the accessed memory address to the *read-set*.

- `WriteTx` is simpler: add or update the location in the *redo-log*.

- `CommitTx` successively
  1. picks up locks for each written object
  2. increments the global version
  3. checks the read objects for being up to date

  before writing redo-log entries to memory while updating their version and realasing their locks

## Properties of TL2

Opacity is guaranteed by aborting on a read accessing an inconsistent value:



Other observations:

- read-only transactions just need to check that read versions are consistent (no need to increment the global clock)
- writing values still requires locks
  - ▶ deadlocks are still possible
  - ▶ since other transactions can be aborted, one can *preempt* transactions that are deadlocked
  - ▶ since lock accesses are generated, computing a lock order up-front might be possible
- there might be contention on the global clock

Executing `atomic` blocks by repeatedly trying to execute them non-atomically creates new problems:

- a transaction might unnecessarily be aborted
  - ▶ the granularity of what is locked might be too large
  - ▶ a TM implementation might impose restrictions:

```
// Thread 1                      // Thread 2
atomic { // clock=12
    ...                          atomic {
                                   WriteTx(&x,0) = 42; // clock=13
                                 }

    int r = ReadTx(&x,0);
} // tx.RV==12 != clock
```

- lock-based commits can cause contention
  - ▶ organize cells that participate in a transaction in one object
  - ▶ compute a new object as result of a transaction
  - ▶ atomically replace a pointer to the old object with a pointer to the new object if the old object has not changed
  - ⤳ idea of the original STM proposal
- TM system should figure out which memory locations must be logged
- danger of live-locks: transaction B might abort A which might abort B ...

Allowing access to other resources than memory inside an `atomic` block poses problems:

- storage management, condition variables, `volatile` variables, input/output
- semantics should be as if `atomic` implements SLA or TSC semantics

# Integrating Non-TM Resources

Allowing access to other resources than memory inside an `atomic` block poses problems:

- storage management, condition variables, `volatile` variables, input/output
- semantics should be as if `atomic` implements SLA or TSC semantics

Usual choice is one of the following:

- *Prohibit It.* Certain constructs do not make sense. Use compiler to reject these programs.
- *Execute It.* I/O operations may only happen in some runs (e.g. file writes usually go to a buffer). Abort if I/O happens.
- *Irrevocably Execute It.* Universal way to deal with operations that cannot be undone: enforce that this transaction terminates (possibly before starting) by making all other transactions conflict.
- *Integrate It.* Re-write code to be transactional: error logging, writing data to a file, . . ..

# Integrating Non-TM Resources

Allowing access to other resources than memory inside an `atomic` block poses problems:

- storage management, condition variables, `volatile` variables, input/output
- semantics should be as if `atomic` implements SLA or TSC semantics

Usual choice is one of the following:

- *Prohibit It.* Certain constructs do not make sense. Use compiler to reject these programs.
- *Execute It.* I/O operations may only happen in some runs (e.g. file writes usually go to a buffer). Abort if I/O happens.
- *Irrevocably Execute It.* Universal way to deal with operations that cannot be undone: enforce that this transaction terminates (possibly before starting) by making all other transactions conflict.
- *Integrate It.* Re-write code to be transactional: error logging, writing data to a file, . . ..

⇝ currently best to use TM only for memory; check if TM supports irrevocable transactions

**Hardware Transactional Memory**

Transactions of a limited size can also be implemented in hardware:

- additional hardware to track read- and write-sets
- conflict detection is *eager* using the cache:
  - additional hardware makes it cheap to perform conflict detection
  - if a cache-line in the read set is invalidated, the transaction aborts
  - if a cache-line in the write set must be written-back, the transaction aborts

⤳ limited by fixed hardware resources, a software backup must be provided

# Hardware Transactional Memory

Transactions of a limited size can also be implemented in hardware:

- additional hardware to track read- and write-sets
- conflict detection is *eager* using the cache:
  - ▶ additional hardware makes it cheap to perform conflict detection
  - ▶ if a cache-line in the read set is invalidated, the transaction aborts
  - ▶ if a cache-line in the write set must be written-back, the transaction aborts

⤳ limited by fixed hardware resources, a software backup must be provided

Two principal implementation of HTM:

1. Explicit Transactional Memory: each access is marked as transactional
   - ▶ similar to `StartTx`, `ReadTx`, `WriteTx`, and `CommitTx`
   - ▶ requires separate transaction instructions
   - ⤳ a transaction has to be translated differently
   - ⚠ mixing transactional and non-transactional accesses is problematic

2. Implicit Transactional Memory: only the beginning and end of a transaction are marked
   - ▶ same instructions can be used, hardware interprets them as transactional
   - ▶ only instructions affecting memory that can be cached can be executed transactionally
   - ▶ hardware access, OS calls, page table changes, etc. all abort a transaction
   - ⤳ provides *strong isolation*

AMD Advanced Synchronization Facilities (ASF):

- defines a logical *speculative region*
- `LOCK MOV` instructions provide *explicit* data transfer between normal memory and speculative region
- aimed to implement larger atomic operations

## Example for HTM

AMD Advanced Synchronization Facilities (ASF):

- defines a logical *speculative region*
- `LOCK MOV` instructions provide *explicit* data transfer between normal memory and speculative region
- aimed to implement larger atomic operations

Intel's TSX in Broadwell/Skylake microarchitecture (since Aug 2014):

- *implicitly transactional*, can use normal instructions within transactions
- tracks read/write set using a single *transaction* bit on cache lines
- provides space for a backup of the whole CPU state (registers, ...)
- use a simple counter to support nested transactions
- may abort at any time due to lack of resources
- aborting in an inner transaction means aborting all of them

## **Example for HTM**

AMD Advanced Synchronization Facilities (ASF):

- defines a logical *speculative region*
- LOCK MOV instructions provide *explicit* data transfer between normal memory and speculative region
- aimed to implement larger atomic operations

Intel's TSX in Broadwell/Skylake microarchitecture (since Aug 2014):

- *implicitly transactional*, can use normal instructions within transactions
- tracks read/write set using a single *transaction* bit on cache lines
- provides space for a backup of the whole CPU state (registers, ...)
- use a simple counter to support nested transactions
- may abort at any time due to lack of resources
- aborting in an inner transaction means aborting all of them

Intel provides two software interfaces to TM:

1. Restricted Transactional Memory (RTM)
2. Hardware Lock Elision (HLE)

**Restricted Transactional Memory**

Supporting Transactional operations:

- augment each cache line with an extra bit *T*
- introduce a nesting counter $C$ and a backup register set

## Implementing RTM using the Cache (Intel)

Supporting Transactional operations:

- augment each cache line with an extra bit *T*
- introduce a nesting counter $C$ and a backup register set



$\rightsquigarrow$ additional transaction logic:

- xbegin increments $C$ and, if $C = 0$, backs up registers and flushes buffer
  - subsequent read or write access to a cache line sets *T* if $C > 0$
  - applying an *invalidate* message to a cache line with *T* flag issues xabort
  - observing a *read* for a *modified* cache line with *T* flag issues xabort
- xabort clears all *T* flags and the store buffer, invalidates the former *TM* lines, sets $C = 0$ and restores CPU registers
- xend decrements $C$ and, if $C = 0$, clears all *T* flags, flushes store buffer

Provides new instructions `xbegin`, `xend`, `xabort`, and `xtest`:

- `xbegin` *on transaction start* skips to the next instruction or *on abort*
  - continues at the given address
  - implicitly stores an error code in `eax`
- `xend` commits the transaction started by the most recent `xbegin`
- `xabort` aborts the whole transaction with an error code
- `xtest` checks if the processor is executing transactionally

# Restricted Transactional Memory

Provides new instructions `xbegin`, `xend`, `xabort`, and `xtest`:

- `xbegin` *on transaction start* skips to the next instruction or *on abort*
  - continues at the given address
  - implicitly stores an error code in `eax`
- `xend` commits the transaction started by the most recent `xbegin`
- `xabort` aborts the whole transaction with an error code
- `xtest` checks if the processor is executing transactionally

The instruction `xbegin` is made accessible via library function `_xbegin()`:

```
_xbegin()

move    eax, 0xFFFFFFFF
xbegin _txnL1
_txnL1:
move    retval, eax
```

# Restricted Transactional Memory

Provides new instructions `xbegin`, `xend`, `xabort`, and `xtest`:

- `xbegin` *on transaction start* skips to the next instruction or *on abort*
  - continues at the given address
  - implicitly stores an error code in `eax`
- `xend` commits the transaction started by the most recent `xbegin`
- `xabort` aborts the whole transaction with an error code
- `xtest` checks if the processor is executing transactionally

The instruction `xbegin` is made accessible via library function `_xbegin()`:

```
_xbegin()

move   eax, 0xFFFFFFFF
xbegin _txnL1
_txnL1:
move   retval, eax
```

```
if(_xbegin()==_XBEGIN_STARTED) {
  // transaction code
  _xend();
} else {
  // non-transactional fall-back
}
```

## Restricted Transactional Memory

Provides new instructions `xbegin`, `xend`, `xabort`, and `xtest`:

- `xbegin` *on transaction start* skips to the next instruction or *on abort*
  - continues at the given address
  - implicitly stores an error code in `eax`
- `xend` commits the transaction started by the most recent `xbegin`
- `xabort` aborts the whole transaction with an error code
- `xtest` checks if the processor is executing transactionally

The instruction `xbegin` is made accessible via library function `_xbegin()`:

```
_xbegin()

move    eax, 0xFFFFFFFF
xbegin _txnL1
_txnL1:
move    retval, eax
```

```
if(_xbegin()==_XBEGIN_STARTED) {
  // transaction code
  _xend();
} else {
  // non-transactional fall-back
}
```

↝ user must provide *fall-back code*

## Considerations for the Fall-Back Path

Consider executing the following code concurrently with itself:

```
int data[100]; // shared
void update(int idx, int value) {
  if(_xbegin()==_XBEGIN_STARTED) {
      data[idx] += value;
      _xend();
    } else {
      data[idx] += value;
    }
}
```

# Considerations for the Fall-Back Path

Consider executing the following code concurrently with itself:

```
int data[100]; // shared
void update(int idx, int value) {
  if(_xbegin()==_XBEGIN_STARTED) {
      data[idx] += value;
      _xend();
    } else {
      data[idx] += value;
    }
}
```

⚠ Several problems:
- the fall-back code may execute racing itself
- the fall-back code is not isolated from the transaction

# Considerations for the Fall-Back Path

Consider executing the following code concurrently with itself:

```
int data[100]; // shared
void update(int idx, int value) {
  if(_xbegin()==_XBEGIN_STARTED) {
      data[idx] += value;
      _xend();
    } else {
      data[idx] += value;
    }
}
```

⚠ Several problems:
- the fall-back code may execute racing itself
- the fall-back code is not isolated from the transaction

⤳ First idea: ensure that the fall-back path is executed atomically

## Protecting the Fall-Back Path

Use a lock to prevent the transaction from interrupting the fall-back path:

```
int data[100]; // shared
int mutex;
void update(int idx, int value) {
  if(_xbegin()==_XBEGIN_STARTED) {

      data[idx] += value;
      _xend();
    } else {
      wait(mutex);
      data[idx] += value;
      signal(mutex);
    }
}
```

- the fall-back code does not execute racing itself ✓

Use a lock to prevent the transaction from interrupting the fall-back path:

```c
int data[100]; // shared
int mutex;
void update(int idx, int value) {
  if(_xbegin()==_XBEGIN_STARTED) {

      data[idx] += value;
      _xend();
    } else {
      wait(mutex);
      data[idx] += value;
      signal(mutex);
    }
}
```

- the fall-back code does not execute racing itself ✓
- ⚠ the fall-back code is still not isolated from the transaction

Use a lock to prevent the transaction from interrupting the fall-back path:

```c
int data[100]; // shared
int mutex;
void update(int idx, int value) {
  if(_xbegin()==_XBEGIN_STARTED) {
      if (!mutex>0) _xabort();
      data[idx] += value;
      _xend();
    } else {
      wait(mutex);
      data[idx] += value;
      signal(mutex);
    }
}
```

- the fall-back code does not execute racing itself ✓
- the fall-back code is now isolated from the transaction ✓

# Happened Before Diagram for Transactions

Augment MESI states with extra bit $T$. CPU A: d:E5 t:E0, CPU B: d:I, tmp/value registers

| Thread A | Thread B |
|---|---|
| ```int t = _xbegin();``` | ```_xbegin();``` |
| ```int tmp = data[idx];``` | ```int tmp = data[idx];``` |
| ```data[idx] = tmp + value;``` | ```data[idx] = tmp + value;``` |
| ```_xend();``` | ```_xend();``` |

## Common Code Pattern for Mutexes

Using HTM in order to implement mutex:

```
int data[100]; // shared
int mutex;
void update(int idx, int val) {
  if(_xbegin()==_XBEGIN_STARTED) {
    if (!mutex>0) _xabort();
    data[idx] += val;
    _xend();
  } else {
    wait(mutex);
    data[idx] += val;
    signal(mutex);
  }
}
```

## Common Code Pattern for Mutexes

Using HTM in order to implement mutex:

```
int data[100]; // shared
int mutex;
void update(int idx, int val) {
  if(_xbegin()==_XBEGIN_STARTED) {
    if (!mutex>0) _xabort();
    data[idx] += val;
    _xend();
  } else {
    wait(mutex);
    data[idx] += val;
    signal(mutex);
  }
}
```

```
void update(int idx, int val) {
  lock(&mutex);
  data[idx] += val;
  unlock(&mutex);
}
void lock(int* mutex) {
  if(_xbegin()==_XBEGIN_STARTED)
  { if (!*mutex>0) _xabort();
    else return;
  } wait(mutex);
}
void unlock(int* mutex) {
  if (!*mutex>0) signal(mutex);
  else _xend();
}
```

- critical section may be executed without taking the lock (the lock is *elided*)
- as soon as one thread conflicts, it aborts, takes the lock in the fallback path and thereby aborts all other transactions that have read `mutex`

**Hardware Lock Elision**

# Hardware Lock Elision

*Observation:* Using RTM to implement lock elision is a common pattern
⤳ provide special handling in hardware: HLE

> **Idea: Hardware Lock Elision**
>
> **①** By default defer actual acquisition of the lock
> **②** Instead rely on HTM to sort out conflicting concurrent accesses
> **③** Fall back to actual locking only in case of conflicts
> **④** Support legacy lock code by locally acting as if semaphore value is actually modified

- requires annotations for lock instructions:
  - ▸ instruction that increments the semaphore must be prefixed with `xacquire`
  - ▸ instruction setting the semaphore to $0$ must be prefixed with `xrelease`
  - ▸ these prefixes are ignored on older platforms
- for a successful elision, all signal/wait operations of a lock must be annotated

Transactional operation:

- re-uses infrastructure for Restricted Transactional Memory
- add a buffer for elided locks, similar to store buffer

## Implementing Lock Elision

Transactional operation:

- re-uses infrastructure for Restricted Transactional Memory
- add a buffer for elided locks, similar to store buffer



- `xacquire` of lock ensures *shared/exclusive* cache line state with *T*, issues `xbegin` and keeps the modified lock value in *elided lock* buffer
  - ▶ r/w access to other cache lines sets *T*
  - ▶ applying an *invalidate* message to a *T* cache line issues `xabort`, analogous for *read* message to a *TM* cache line
  - ▶ a *local CPU load* from the address of the elided lock accesses the buffer
- on `xrelease` on the same lock, decrement $C$ and, if $C = 0$, clear *T* flags and elided locks buffer flush the store buffer

# Transactional Memory: Summary

Transactional memory aims to provide `atomic` blocks for general code:

- frees the user from deciding how to lock data structures
- compositional way of communicating concurrently
- can be implemented using software (locks, atomic updates) or hardware

# Transactional Memory: Summary

Transactional memory aims to provide `atomic` blocks for general code:

- frees the user from deciding how to lock data structures
- compositional way of communicating concurrently
- can be implemented using software (locks, atomic updates) or hardware

It is hard to get the details right:

- semantics of *explicit HTM* and *STM* transactions quite subtle when mixing with non-TM (*weak* vs. *strong isolation*)
- *single-lock atomicity* vs. *transactional sequential consistency* semantics
- STM not the right tool to synchronize threads without shared variables
- TM providing *opacity* (serializability) requires *eager conflict detection* or *lazy version management*

Pitfalls in *implicit* HTM:

- RTM requires a fall-back path
- no progress guarantee
- HLE can be implemented in software using RTM

# TM in Practice

Availability of TM Implementations:

- GCC can translate accesses in `__transaction_atomic` regions into `libitm` library calls
- the library `libitm` provides different TM implementations:
  1. On systems with TSX, it maps atomic blocks to HTM instructions
  2. On systems without TSX and for the fallback path, it resorts to STM
- C++20 standardizes `synchronized`/`atomic_XXX` blocks
- RTM support slowly introduced to OpenJDK Hotspot monitors

Availability of TM Implementations:

- GCC can translate accesses in `__transaction_atomic` regions into `libitm` library calls
- the library `libitm` provides different TM implementations:
  1. On systems with TSX, it maps atomic blocks to HTM instructions
  2. On systems without TSX and for the fallback path, it resorts to STM
- C++20 standardizes `synchronized`/`atomic_XXX` blocks
- RTM support slowly introduced to OpenJDK Hotspot monitors

Use of hardware lock elision is limited:

- allows to easily convert existing locks
- `pthread` locks in `glibc` use RTM https://lwn.net/Articles/534758/:
  - allows implementation of fallback mechanisms
  - HLE only special case of general lock
- implementing monitors is challenging
  - lock count and thread id may lead to conflicting accesses
  - in `pthreads`: error conditions often not checked anymore

# **Outlook**

Several other principles exist for concurrent programming:

1. non-blocking message passing (the actor model)
   - a program consists of actors that send messages
   - each actor has a queue of incoming messages
   - messages can be processed and new messages can be sent
   - special filtering of incoming messages
   - *example:* Erlang, many add-ons to existing languages

2. blocking message passing (CSP, $\pi$-calculus, join-calculus)
   - a process sends a message over a channel and blocks until the recipient accepts it
   - channels can be send over channels ($\pi$-calculus)
   - *examples:* Occam, Occam-$\pi$, Go

3. (immediate) priority ceiling
   - declare *processes* with priority and *resources* that each process may acquire
   - each resource has the maximum (ceiling) priority of all processes that may acquire it
   - a process' priority at run-time increases to the maximum of the priorities of held resources
   - the process with the maximum (run-time) priority executes

# References

TUΠ

📕 D. Dice, O. Shalev, and N. Shavit.
Transactional Locking II.
In *Distributed Coputing*, LNCS, pages 194–208. Springer, Sept. 2006.

📕 T. Harris, J. Larus, and R. Rajwar.
Transactional memory, 2nd edition.
*Synthesis Lectures on Computer Architecture*, 5(1):1–263, 2010.

Online resources on Intel HTM and GCC's STM:

1. http://software.intel.com/en-us/blogs/2013/07/25/
   fun-with-intel-transactional-synchronization-extensions

2. http://www.realworldtech.com/haswell-tm/4/

3. http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2012/n3341.pdf

# **Programming Languages**

Dispatching Method Calls

Dr. Michael Petter
Winter Term 2019

# Dispatching - Outline

## Dispatching

1. Motivation
2. Formal Model
3. Quiz
4. Dispatching from the Inside

## Solutions in Single-Dispatching

1. Type introspection
2. Generic interface

## Multi-Dispatching

1. Formal Model
2. Multi-Java
3. Multi-dispatching in Perl6
4. Multi-dispatching in Clojure

Section 1

**Direct Function Calls**

# Function Dispatching (ANSI C89)

```c
#include <stdio.h>

void fun(int i) {  }
void bar(int i, double j) {  }

int main(){
  fun(1);
  bar(1,1.2);
  void (*foo)(int);
  foo = fun;
  return 0;
}
```

Section 2

**Overloading Function Names**

# Function Dispatching (ANSI C89)

```c
#include <stdio.h>

void println(int i)   { print("%d\n",i); };
void println(float f) { print("%f\n",f); };

int main(){
  println(1.2);
  println(1);
  return 0;
}
```

```
#include <stdio.h>

void println(int i)   { print("%d\n",i); };
void println(float f) { print("%f\n",f); };

int main(){
  println(1.2);
  println(1);
  return 0;
}
```

⚠ Functions with same names but different parameters not legal

## Generic Selection (C11)

$$\textit{generic-selection} \;\mapsto\; \texttt{\_Generic}(\textit{exp}, \textit{generic-assoclist})$$
$$\textit{generic-assoclist} \;\mapsto\; (\textit{generic-assoc}, )^* \textit{generic-assoc}$$
$$\textit{generic-assoc} \;\mapsto\; \texttt{typename} : \textit{exp} \mid \texttt{default} : \textit{exp}$$

### Example:
```c
#include <stdio.h>
#define printf_dec_format(x) _Generic((x), \
    signed int: "%d", \
    float: "%f" )
#define println(x) printf(printf_dec_format(x), x), printf("\n");

int main(){
  println(1.2);
  println(1);
  return 0;
}
```

## Generic Selection (C11)

$$generic\text{-}selection \mapsto \texttt{\_Generic}(exp, generic\text{-}assoclist)$$
$$generic\text{-}assoclist \mapsto (generic\text{-}assoc,)^* generic\text{-}assoc$$
$$generic\text{-}assoc \mapsto \texttt{typename} : exp \mid \texttt{default} : exp$$

Example:
```c
#include <stdio.h>

int main(){
  printf(_Generic((1.2),signed int: "%d",float: "%f"), 1.2), printf("\n");
  printf(_Generic((  1),signed int: "%d",float: "%f"),   1), printf("\n");
  return 0;
}
```

```
class D {
  public static void p(Object o) { System.out.print(o); }
  public         int f(int i)    { p("f(int): ");    return i+1; }
  public      double f(double d) { p("f(double): "); return d+1.3;}
}

public static void main() {
  D d = new D();
  D.p(d.f(2)+"\n");
  D.p(d.f(2.3)+"\n");
}
```

# Overloading (Java/C++)

```java
class D {
  public static void p(Object o) { System.out.print(o); }
  public         int f(int i)    { p("f(int): ");    return i+1; }
  public      double f(double d) { p("f(double): "); return d+1.3;}
}

public static void main() {
  D d = new D();
  D.p(d.f(2)+"\n");
  D.p(d.f(2.3)+"\n");
}
```

```
>$ javac Overloading.java; java Overloading
f(int): 3
f(double): 3.6
```

## Overloading with Inheritance (Java)

```java
class B {
  public static void p(Object o) { System.out.print(o); }
  public        int f(int i)    { p("f(int): ");    return i+1; }
}
class D extends B {
  public     double f(double d) { p("f(double): "); return d+1.3;}
}

public static void main() {
  D d = new D();
  B.p(d.f(2)+"\n");
  B.p(d.f(2.3)+"\n");
}
```

## Overloading with Inheritance (Java)

```java
class B {
  public static void p(Object o) { System.out.print(o); }
  public        int f(int i)     { p("f(int): ");     return i+1; }
}
class D extends B {
  public     double f(double d) { p("f(double): "); return d+1.3;}
}

public static void main() {
  D d = new D();
  B.p(d.f(2)+"\n");
  B.p(d.f(2.3)+"\n");
}
```

```
>$ javac Overloading.java; java Overloading
f(int): 3
f(double): 3.6
```

## Overloading with Scopes (C++)

```cpp
#include<iostream>
using namespace std;
class B { public:
  int f(int i) { cout << "f(int): "; return i+1; }
};
class D : public B { public:

  double f(double d) { cout << "f(double): "; return d+1.3; }
};

int main() {
  D* pd = new D;
  cout << pd->f(2) << '\n';
  cout << pd->f(2.3) << '\n';
}
```

## Overloading with Scopes (C++)

```cpp
#include<iostream>
using namespace std;
class B { public:
  int f(int i) { cout << "f(int): "; return i+1; }
};
class D : public B { public:

  double f(double d) { cout << "f(double): "; return d+1.3; }
};

int main() {
  D* pd = new D;
  cout << pd->f(2) << '\n';
  cout << pd->f(2.3) << '\n';
}
```

```
>$ ./overloading
f(double): 3.3
f(double): 3.6
```

## Overloading with Scopes (C++)

```cpp
#include<iostream>
using namespace std;
class B { public:
  int f(int i) { cout << "f(int): "; return i+1; }
};
class D : public B { public:
  using B::f;
  double f(double d) { cout << "f(double): "; return d+1.3; }
};

int main() {
  D* pd = new D;
  cout << pd->f(2) << '\n';
  cout << pd->f(2.3) << '\n';
}
```

```
>$ ./overloading
f(int): 3
f(double): 3.6
```

```
class D {
  public static void p(Object o) { System.out.print(o); }
  public          int f(int i, double j) { p("f(i,d): "); return i;}
  public          int f(double i, int j) { p("f(d,i): "); return j;}
}

public static void main() {
  D d = new D();
  D.p(d.f(2,2)+"\n");
}
```

```
class D {
  public static void p(Object o) { System.out.print(o); }
  public          int f(int i, double j) { p("f(i,d): "); return i;}
  public          int f(double i, int j) { p("f(d,i): "); return j;}
}

public static void main() {
  D d = new D();
  D.p(d.f(2,2)+"\n");
}
```

⚠

```
>$ javac Overloading.java
Overloading.java:(?): error: reference to f is ambiguous
```

# Static Methods are *Statically Dispatched*

# Static Methods are *Statically Dispatched*

**Function Call Expression**

Function to be dispatched

# Static Methods are *Statically Dispatched*

**Function Call Expression**

Function to be dispatched



**Signature**

- Function Name
- *Static* Types of Parameters
- Return Type

# Static Methods are *Statically Dispatched*

**Function Call Expression**

Function to be dispatched

**Concrete Method**

Provides calling target for a call signature

$f'(e_1,...,e_n)$ — dispatches to / handles → $t_0\ f(t_1\ p_1,...,t_n\ p_n)$

determines ↓

is applicable to ↗

**Signature**
$t_0',..., t_n'$

**Signature**

- Function Name
- *Static* Types of Parameters
- Return Type

# Static Methods are *Statically Dispatched*

**Function Call Expression**

Function to be dispatched

**Concrete Method**

Provides calling target for a call signature



$f'(e_1,...,e_n)$ — dispatches to / handles → $t_0\ f(t_1\ p_1,...,t_n\ p_n)$

determines → **Signature** $t_0',..., t_n'$

is applicable to

**Signature**

- Function Name
- *Static* Types of Parameters
- Return Type

$f$ **is applicable to** $f' \Leftrightarrow f \leq f'$**:**

$\leq$ is the *subtype relation*:

$$R\, f(T_1, \ldots, T_n) \leq R'\, f'(T_1', \ldots, T_n')$$

$$\implies R \leq R' \wedge T_i' \leq T_i$$

## Inside the Javac – Predicates

Concept of methods being *applicable* for arguments:

```java
// true if the given method is applicable to the given arguments
boolean isApplicable(MemberDefinition m, Type args[]) {
    // Sanity checks:
    Type mType = m.getType();
    if (!mType.isType(TC_METHOD))            return false;

    Type mArgs[] = mType.getArgumentTypes();
    if (args.length != mArgs.length)         return false;

    for (int i = args.length ; --i >= 0 ;)
        if (!isMoreSpecific(args[i], mArgs[i])) return false;
    return true;
}
boolean isMoreSpecific(Type moreSpec, Type lessSpec) //... type based specialization
```

Concept of method signatures being *more specific* then others:

```java
// true if "more" is in every argument at least as specific as "less"
boolean isMoreSpecific(MemberDefinition more, MemberDefinition less) {
    Type moreType = more.getClassDeclaration().getType();
    Type lessType = less.getClassDeclaration().getType();
    return isMoreSpecific(moreType, lessType) // return type based comparison
        && isApplicable(less, more.getType().getArgumentTypes()); // parameter type based
}
```

# Finding the Most Specific Concrete Method

```
MemberDefinition matchMethod(Environment env, ClassDefinition accessor,
                             Identifier methodName, Type[] argumentTypes) throws ... {
    // A tentative maximally specific method.
    MemberDefinition tentative = null;
    // A list of other methods which may be maximally specific too.
    List candidateList = null;
    // Get all the methods inherited by this class which have the name `methodName'
    for (MemberDefinition method : allMethods.lookupName(methodName)) {
        // See if this method is applicable.
        if (!env.isApplicable(method, argumentTypes)) continue;
        // See if this method is accessible.
        if ((accessor != null) && (!accessor.canAccess(env, method))) continue;
        if ((tentative == null) || (env.isMoreSpecific(method, tentative)))
            // `method' becomes our tentative maximally specific match.
            tentative = method;
        else {  // If this method could possibly be another maximally specific
                // method, add it to our list of other candidates.
                if (!env.isMoreSpecific(tentative,method)) {
                    if (candidateList == null) candidateList = new ArrayList();
                    candidateList.add(method);
    }   }      }
    if (tentative != null && candidateList != null)
        // Find out if our `tentative' match is a uniquely maximally specific.
        for (MemberDefinition method : candidateList )
            if (!env.isMoreSpecific(tentative, method))
                throw new AmbiguousMember(tentative, method);
    return tentative;
}
```

Section 3

**Overriding Methods**

## Object Orientation

**Emphasizing the *Receiver* of a Call**

In Object Orientation, we see objects associating strongly with particular procedures, a.k.a. *Methods*.

```java
class Natural {
  int value;
}
void incBy(Natural n,int i) {
  n.value += Math.abs(i);
}

...
incBy(nat,42);
```

$\Longrightarrow$

```java
class Natural {
  int value;

  void incBy(int i){
    this.value += Math.abs(i);
  }
}
...
nat.incBy(42);
```

- Associating the first parameter as *Receiver* of the method, and pulling it out of the parameters list
- Implicitely binding the first parameter to the fixed name *this*

# Subtyping in Object Orientation

**Emphasizing the *Receiver*'s Responsibility**

An Object Oriented Subtype is supposed to take responsibility for calls to Methods that are associated with the type, that it specializes.

```java
class Integral {
  int i;
  void incBy(int delta){
      i += delta;
  }
}
class Natural extends Integral {
  int value;
  void incBy(int i){
    this.value += Math.abs(i);
  }
}
```

```java
Integral i = new Integral(-5);
i.incBy(42);
Natural n = new Natural(42);
n.incBy(42);
i = n;
i.incBy(42);
```

⚠ In OO, at runtime subtypes can inhabit statically more general typed variables

⤳ Implicitly call the specialized method!

# Methods are *dynamically dispatched*

**Function Call Expression**

Call expression to be dispatched.

**Concrete Method**

Provides calling target for a call signature

$f'(e_1,...,e_n)$ — dispatches to / handles → $t_0 \; f(t_1 \; p_1,...,t_n \; p_n)$

determines ↓

is applicable to ↗

**Signature**
$t_0',..., t_n'$

**Signature**

Static types of actual parameters.

# Methods are *dynamically dispatched*

**Function Call Expression**

Call expression to be dispatched.

**Concrete Method**

Provides calling target for a call signature

$f'(e_1,...,e_n)$ — dispatches to / handles — $t_0\ f(t_1\ p_1,...,t_n\ p_n)$

determines → **Signature** $t_0',...,t_n'$

is applicable to

specialized by → **Specializer** $t_0^*, t_1^*$

**Signature**

Static types of actual parameters.

**Specializer**

Specialized types to be matched at the call

# How can we implement that?

**Let's look at what Java does!**

The Java platform as example for state of the art OO systems:

- Static Javac-based compiler
- Dynamic Hotspot JIT-Compiler/Interpreter

Let's watch the following code on its way to the CPU:

```java
public static void main(String[] args){
    Integral i = new Natural(1);
    i.incBy(42);
}
```

$\rightsquigarrow$ `matchMethod` returns the statically most specific signature

$\rightsquigarrow$ Codegeneration hardcodes `invokevirtual` with this signature

```
Code:
   0: new             #4                  // class Natural
   3: dup
   4: iconst_1
   5: invokespecial   #5                  // Method "<init>":(I)V
   8: astore_1
   9: aload_1
  10: bipush          42
  12: invokevirtual   #6                  // Method Integral.incBy:(I)V
  15: return
```

**?** What is the semantics of `invokevirtual`?

$\rightsquigarrow$ `matchMethod` returns the statically most specific signature

$\rightsquigarrow$ Codegeneration hardcodes `invokevirtual` with this signature

```
Code:
   0: new          #4              // class Natural
   3: dup
   4: iconst_1
   5: invokespecial #5             // Method "<init>":(I)V
   8: astore_1
   9: aload_1
  10: bipush       42
  12: invokevirtual #6             // Method Integral.incBy:(I)V
  15: return
```

**?** What is the semantics of `invokevirtual`?

$\rightsquigarrow$ Check the runtime interpreter: Hotspot VM calls `resolve_method`!

TIM

```
void LinkResolver::resolve_method(methodHandle& resolved_method, KlassHandle resolved_klass,
                                  Symbol* method_name, Symbol* method_signature,
                                  KlassHandle current_klass) {

  // 1. check if klass is not interface
  if (resolved_klass->is_interface()) ;//... throw "Found interface, but class was expected"

  // 2. lookup method in resolved klass and its super klasses
  lookup_method_in_klasses(resolved_method, resolved_klass, method_name, method_signature);
      // calls klass::lookup_method() -> next slide

  if (resolved_method.is_null()) { // not found in the class hierarchy
    // 3. lookup method in all the interfaces implemented by the resolved klass
    lookup_method_in_interfaces(resolved_method, resolved_klass, method_name, method_signature);

    if (resolved_method.is_null()) {
      // JSR 292: see if this is an implicitly generated method MethodHandle.invoke(*...)
      lookup_implicit_method(resolved_method, resolved_klass, method_name, method_signature, current_klass);
    }

    if (resolved_method.is_null()) { // 4. method lookup failed
      // ... throw  java_lang_NoSuchMethodError()
  } }

  // 5. check if method is concrete
  if (resolved_method->is_abstract() && !resolved_klass->is_abstract()) {
    // ... throw java_lang_AbstractMethodError()
  }

  // 6. access checks, etc.
}
```

The method lookup recursively traverses the super class chain:

```
MethodDesc* klass::lookup_method(Symbol* name, Symbol* signature) {
  for (KlassDesc* klas = as_klassOop(); klas != NULL; klas = klass::cast(klas)->super()) {
    MethodDesc* method = klass::cast(klass)->find_method(name, signature);
    if (method != NULL) return method;
  }
  return NULL;
}
```

# Inside the Hotspot VM

```
MethodDesc* klass::find_method(ObjArrayDesc* methods, Symbol* name, Symbol* signature) {
  int len = methods->length();
  // methods are sorted, so do binary search
  int i, l = 0 , h = len - 1;
  while (l <= h) {
    int mid = (l + h) >> 1;
    MethodDesc* m = (MethodDesc*)methods->obj_at(mid);
    int res = m->name()->fast_compare(name);
    if (res == 0) {
      // found matching name; do linear search to find matching signature
      // first, quick check for common case
      if (m->signature() == signature) return m;
      // search downwards through overloaded methods
      for (i = mid - 1; i >= l; i--) {
        MethodDesc* m = (MethodDesc*)methods->obj_at(i);
        if (m->name() != name) break;
        if (m->signature() == signature) return m;
      }
      // search upwards
      for (i = mid + 1; i <= h; i++) {
        MethodDesc* m = (MethodDesc*)methods->obj_at(i);
        if (m->name() != name) break;
        if (m->signature() == signature) return m;
      }
      return NULL; // not found
    } else if (res < 0) l = mid + 1;
           else           h = mid - 1;
  }
  return NULL;
}
```

# Single-Dispatching: Summary



**Compile Time**

**Runtime**

### Javac

Matches a method call expression *statically* to the *most specific* method signature via `matchMethod( ... )`

### Hotspot VM

Interprets `invokevirtual` via `resolve_method(...)`, scanning the superclass chain with `find_method(...)` for the statically fixed `signature`

## Example: Sets of Natural Numbers

```java
class Natural {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean equals(Natural n){
    return n.number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

## Example: Sets of Natural Numbers

```java
class Natural {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean equals(Natural n){
    return n.number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

```
>$ java Natural
[0,0]
```

⚠ Why? Is HashSet buggy?

## Example: Sets of Natural Numbers

```java
class Natural {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean equals(Natural n){
    return n.number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

```
>$ java Natural
[0,0]
```

⚠ Why? Is HashSet buggy?
⤳ Keep attention to exact signature!

## Mini-Quiz: Java Method Dispatching

```java
class A {
  public static  void p  (Object o) { System.out.println(o); }
  public         void m1 (A a) {       p("m1(A) in A");        }
  public         void m1 () {          m1(new B());            }
  public         void m2 (A a) {       p("m2(A) in A");        }
  public         void m2 () {          m2(this);               }
}
class B extends A {
  public         void m1 (B b) {       p("m1(B) in B");        }
  public         void m2 (A a) {       p("m2(A) in B");        }
  public         void m3 () {          super.m1(this);         }
}
```

```java
B b = new B(); A a = b; a.m1(b);
```

## Mini-Quiz: Java Method Dispatching

```java
class A {
  public static  void p  (Object o) { System.out.println(o); }
  public         void m1 (A a) {       p("m1(A) in A");      }
  public         void m1 () {          m1(new B());          }
  public         void m2 (A a) {       p("m2(A) in A");      }
  public         void m2 () {          m2(this);             }
}
class B extends A {
  public         void m1 (B b) {       p("m1(B) in B");      }
  public         void m2 (A a) {       p("m2(A) in B");      }
  public         void m3 () {          super.m1(this);       }
}
```

```
B b = new B(); A a = b; a.m1(b);                                 m1(A) in A
B b = new B(); B a = b; b.m1(a);
```

## Mini-Quiz: Java Method Dispatching

```java
class A {
  public static  void p  (Object o) { System.out.println(o); }
  public         void m1 (A a) {       p("m1(A) in A");        }
  public         void m1 () {          m1(new B());            }
  public         void m2 (A a) {       p("m2(A) in A");        }
  public         void m2 () {          m2(this);               }
}
class B extends A {
  public         void m1 (B b) {       p("m1(B) in B");        }
  public         void m2 (A a) {       p("m2(A) in B");        }
  public         void m3 () {          super.m1(this);         }
}
```

```java
B b = new B(); A a = b; a.m1(b);
B b = new B(); B a = b; b.m1(a);
B b = new B(); b.m2();
```

```
m1(A) in A
m1(B) in B
```

## Mini-Quiz: Java Method Dispatching

```java
class A {
  public static  void p  (Object o) { System.out.println(o); }
  public         void m1 (A a) {       p("m1(A) in A");       }
  public         void m1 () {          m1(new B());           }
  public         void m2 (A a) {       p("m2(A) in A");       }
  public         void m2 () {          m2(this);              }
}
class B extends A {
  public         void m1 (B b) {       p("m1(B) in B");       }
  public         void m2 (A a) {       p("m2(A) in B");       }
  public         void m3 () {          super.m1(this);        }
}
```

```
B b = new B(); A a = b; a.m1(b);                              m1(A) in A
B b = new B(); B a = b; b.m1(a);                              m1(B) in B
B b = new B(); b.m2();                                        m2(A) in B
B b = new B(); b.m1();
```

## Mini-Quiz: Java Method Dispatching

```java
class A {
  public static  void p  (Object o) { System.out.println(o); }
  public         void m1 (A a) {      p("m1(A) in A");      }
  public         void m1 () {         m1(new B());          }
  public         void m2 (A a) {      p("m2(A) in A");      }
  public         void m2 () {         m2(this);             }
}
class B extends A {
  public         void m1 (B b) {      p("m1(B) in B");      }
  public         void m2 (A a) {      p("m2(A) in B");      }
  public         void m3 () {         super.m1(this);       }
}
```

```
B b = new B(); A a = b; a.m1(b);                          m1(A) in A
B b = new B(); B a = b; b.m1(a);                          m1(B) in B
B b = new B(); b.m2();                                    m2(A) in B
B b = new B(); b.m1();                                    m1(A) in A
B b = new B(); b.m3();
```

## Mini-Quiz: Java Method Dispatching

```java
class A {
  public static  void p  (Object o) { System.out.println(o); }
  public         void m1 (A a) {      p("m1(A) in A");        }
  public         void m1 () {         m1(new B());            }
  public         void m2 (A a) {      p("m2(A) in A");        }
  public         void m2 () {         m2(this);               }
}
class B extends A {
  public         void m1 (B b) {      p("m1(B) in B");        }
  public         void m2 (A a) {      p("m2(A) in B");        }
  public         void m3 () {         super.m1(this);         }
}
```

```
B b = new B(); A a = b; a.m1(b);                          m1(A) in A
B b = new B(); B a = b; b.m1(a);                          m1(B) in B
B b = new B(); b.m2();                                    m2(A) in B
B b = new B(); b.m1();                                    m1(A) in A
B b = new B(); b.m3();                                    m1(A) in A
```

Section 4

**Multi-Dispatching**

Mainstream languages support specialization of first parameter:

C++, Java, C#, Smalltalk, Lisp

**So how do we solve the `equals()` problem?**

1. introspection?
2. generic programming?
3. double dispatching?

# Introspection

```
class Natural {
  Natural(int n) { number=Math.abs(n); }
  int number;
  public boolean equals(Object n){
    if (!(n instanceof Natural)) return false;
    return ((Natural)n).number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

```java
class Natural {
  Natural(int n) { number=Math.abs(n); }
  int number;
  public boolean equals(Object n){
    if (!(n instanceof Natural)) return false;
    return ((Natural)n).number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

```
>$ java Natural
[0]
```

✓ Works

# Introspection

```
class Natural {
  Natural(int n) { number=Math.abs(n); }
  int number;
  public boolean equals(Object n){
    if (!(n instanceof Natural)) return false;
    return ((Natural)n).number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

```
>$ java Natural
[0]
```

✓ Works ⚠ but burdens programmer with type safety

## Introspection

```
class Natural {
  Natural(int n) { number=Math.abs(n); }
  int number;
  public boolean equals(Object n){
    if (!(n instanceof Natural)) return false;
    return ((Natural)n).number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

```
>$ java Natural
[0]
```

✓ Works ⚠ but burdens programmer with type safety
⚠ and is only available for languages with type introspection

# Generic Programming

```
interface Equalizable<T>{
  boolean equals(T other);
}
class Natural implements Equalizable<Natural> {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean equals(Natural n){
    return n.number == number;
  }
}
...
  EqualizableAwareSet<Natural> set = new MyHashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

# Generic Programming

```
interface Equalizable<T>{
  boolean equals(T other);
}
class Natural implements Equalizable<Natural> {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean equals(Natural n){
    return n.number == number;
  }
}
...
  EqualizableAwareSet<Natural> set = new MyHashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

⚠ needs another Set implementation and...

# Generic Programming

```java
interface Equalizable<T>{
  boolean equals(T other);
}
class Natural implements Equalizable<Natural> {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean equals(Natural n){
    return n.number == number;
  }
}
...
  EqualizableAwareSet<Natural> set = new MyHashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

⚠ needs another Set implementation and...
⚠ only works for one overloaded version in super hierarchy

```
>$ javac Natural.java
Natural.java:2: error: name clash: equals(T) in Equalizable and equals(Object)
in Object have the same erasure, yet neither overrides the other
```

# Double Dispatching

```
abstract class EqualsDispatcher{
  boolean dispatch(Natural) { return false };
  boolean dispatch(Object)  { return false };
}
class Natural {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean doubleDispatch(EqualsDispatcher ed) {
    return ed.dispatch(this);
  }
  public boolean equals(Object n){
    return n.doubleDispatch(
      new EqualsDispatcher(){
        boolean dispatch(Natural nat) {
          return nat.number==number;
      }; }
); } }
```

# Double Dispatching

```
abstract class EqualsDispatcher{
  boolean dispatch(Natural) { return false };
  boolean dispatch(Object)  { return false };
}
class Natural {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean doubleDispatch(EqualsDispatcher ed) {
    return ed.dispatch(this);
  }
  public boolean equals(Object n){
    return n.doubleDispatch(
        new EqualsDispatcher(){
          boolean dispatch(Natural nat) {
            return nat.number==number;
        }; }
); } }
```

✓ Works

# Double Dispatching

```
abstract class EqualsDispatcher{
  boolean dispatch(Natural) { return false };
  boolean dispatch(Object)  { return false };
}
class Natural {
  Natural(int n){ number=Math.abs(n); }
  int number;
  public boolean doubleDispatch(EqualsDispatcher ed) {
    return ed.dispatch(this);
  }
  public boolean equals(Object n){
    return n.doubleDispatch(
      new EqualsDispatcher(){
        boolean dispatch(Natural nat) {
          return nat.number==number;
      }; }
); } }
```

✓ Works ⚠ but needs Dispatcher to know complete class hierarchies

# Formal Model of Multi-Dispatching [7]

**Idea**

Introduce Specializers for all parameters

$f'(e_1, \ldots, e_n)$ — dispatches to / handles — $t_0 \ f(t_1 \ p_1, \ldots, t_n \ p_n)$

determines

is applicable to

specialized by

**Signature**
$t_0', \ldots, t_n'$

**Specializer**
$t^*_0, t^*_1 \ldots t^*_n$

# Formal Model of Multi-Dispatching [7]



**Idea**

Introduce Specializers for all parameters

In the diagram:

- $f'(e_1,...,e_n)$ — **dispatches to** / **handles** — $t_0\ f(t_1\ p_1,...,t_n\ p_n)$
- $f'(e_1,...,e_n)$ — **determines** — **Signature** $t_0',..., t_n'$
- $t_0\ f(t_1\ p_1,...,t_n\ p_n)$ — **is applicable to** — **Signature** $t_0',..., t_n'$
- $t_0\ f(t_1\ p_1,...,t_n\ p_n)$ — **specialized by** — **Specializer** $t^*_0, t^*_1 ... t^*_n$

**How it works**

1. Specializers as subtype annotations to parameter types
2. Dispatcher selects *Most Specific* Concrete Method

# Implications of the implementation

## Type-Checking

1. Typechecking families of concrete methods introduces checking the existance of unique most specific methods for all *valid visible type tuples*.
2. Multiple-Inheritance or interfaces as specializers introduce ambiguities, and thus induce runtime ambiguity exceptions

## Code-Generation

1. Specialized methods generated separately
2. Dispatcher method calls specialized methods
3. Order of the dispatch tests determines the most specialized method

## Performance penalty

The runtime-penalty for multi-dispatching is related to the number of parameters of a multi-method many `instanceof` tests.

```
class Natural {
  public Natural(int n){ number=Math.abs(n); }
  private int number;
  public boolean equals(Object@Natural n){
    return n.number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

```java
class Natural {
  public Natural(int n){ number=Math.abs(n); }
  private int number;
  public boolean equals(Object@Natural n){
    return n.number == number;
  }
}
...
  Set<Natural> set = new HashSet<>();
  set.add(new Natural(0));
  set.add(new Natural(0));
  System.out.println(set);
```

```
>$ java Natural
[0]
```

✓ Clean Code!

## Natural Numbers Behind the Scenes

```
>$ javap -c Natural
```

```
public boolean equals(java.lang.Object);
 Code:
  0:   aload_1
  1:   instanceof      #2; //class Natural
  4:   ifeq    16
  7:   aload_0
  8:   aload_1
  9:   checkcast       #2; //class Natural
  12:  invokespecial   #28; //Method equals$body3$0:(LNatural;)Z
  15:  ireturn
  16:  aload_0
  17:  aload_1
  18:  invokespecial   #31; //Method equals$body3$1:(LObject;)Z
  21:  ireturn
```

## Natural Numbers Behind the Scenes

```
>$ javap -c Natural
```

```
public boolean equals(java.lang.Object);
 Code:
  0:   aload_1
  1:   instanceof      #2; //class Natural
  4:   ifeq    16
  7:   aload_0
  8:   aload_1
  9:   checkcast       #2; //class Natural
  12:  invokespecial   #28; //Method equals$body3$0:(LNatural;)Z
  15:  ireturn
  16:  aload_0
  17:  aload_1
  18:  invokespecial   #31; //Method equals$body3$1:(LObject;)Z
  21:  ireturn
```

⤳ Redirection to methods equals$body3$1 and equals$body3$0

Section 5

**Natively multidispatching Languages**

## Perl6

TUM

```
my Cool $foo;
my Cool $bar;
multi fun(Cool $one, Cool $two){
    say "Dispatch base"
}
multi fun(Int $one,Str $two){
    say "Dispatch 1"
}
multi fun(Str $one,Int $two){
    say "Dispatch 2"
}
$foo=1;
$bar="blabla";
fun($foo,$bar);
```

# Perl6

```perl
my Cool $foo;
my Cool $bar;
multi fun(Cool $one, Cool $two){
    say "Dispatch base"
}
multi fun(Int $one,Str $two){
    say "Dispatch 1"
}
multi fun(Str $one,Int $two){
    say "Dispatch 2"
}
$foo=1;
$bar="blabla";
fun($foo,$bar);
```

```
Dispatch 1
```

# Perl6

```perl6
my Cool $foo;
my Cool $bar;
multi fun(Cool $one, Cool $two){
    say "Dispatch base"
}
multi fun(Int $one,Str $two){
    say "Dispatch 1"
}
multi fun(Str $one,Int $two){
    say "Dispatch 2"
}
$foo=1;
$bar="blabla";
fun($foo,$bar);

$foo="bla";
fun($foo,$bar)
```

```
Dispatch 1
```

```
my Cool $foo;
my Cool $bar;
multi fun(Cool $one, Cool $two){
    say "Dispatch base"
}
multi fun(Int $one,Str $two){
    say "Dispatch 1"
}
multi fun(Str $one,Int $two){
    say "Dispatch 2"
}
$foo=1;
$bar="blabla";
fun($foo,$bar);

$foo="bla";
fun($foo,$bar)
```

```
Dispatch 1
Dispatch base
```

# Clojure

... is a *lisp* dialect for the JVM with:

- Prefix notation
- () – Brackets for lists
- :: – Userdefined keyword constructor ::keyword
- [] – Vector constructor
- `fn` – Creates a lambda expression
  `(fn [x y] (+ x y))`
- `derive` – Generates hierarchical relationships
  `(derive ::child ::parent)`
- `defmulti` – Creates new generic method
  `(defmulti name dispatch-fn)`
- `defmethod` – Creates new concrete method
  `(defmethod name dispatch-val &fn-tail)`

# Principle of Multidispatching in Clojure

```clojure
(derive ::child ::parent)

(defmulti fun (fn [a b] [a b]))
(defmethod fun [::child  ::child ] [a b] "child  equals")
(defmethod fun [::parent ::parent] [a b] "parent equals")

(pr (fun ::child ::child))
```

# Principle of Multidispatching in Clojure

```clojure
(derive ::child ::parent)

(defmulti fun (fn [a b] [a b]))
(defmethod fun [::child  ::child ] [a b] "child  equals")
(defmethod fun [::parent ::parent] [a b] "parent equals")

(pr (fun ::child ::child))
```

```
child equals
```

# More Creative dispatching in Clojure

```clojure
(defn salary [amount]
        (cond (< amount 600)   ::poor
              (>= amount 5000) ::rich
              :else            ::average))

(defrecord UniPerson [name wage])

(defmulti  print (fn [person] (salary (:wage person)) ))
(defmethod print ::poor    [person](str "HiWi  " (:name person)))
(defmethod print ::average [person](str "Dr.   " (:name person)))
(defmethod print ::rich    [person](str "Prof. " (:name person)))

(pr (print (UniPerson. "Petter" 2000)))
(pr (print (UniPerson. "Stefan"  200)))
(pr (print (UniPerson. "Seidl" 16000)))
```

## More Creative dispatching in Clojure

```clojure
(defn salary [amount]
        (cond (< amount 600)   ::poor
              (>= amount 5000) ::rich
              :else            ::average))

(defrecord UniPerson [name wage])

(defmulti  print (fn [person] (salary (:wage person)) ))
(defmethod print ::poor    [person](str "HiWi " (:name person)))
(defmethod print ::average [person](str "Dr.  " (:name person)))
(defmethod print ::rich    [person](str "Prof. " (:name person)))

(pr (print (UniPerson. "Petter" 2000)))
(pr (print (UniPerson. "Stefan"  200)))
(pr (print (UniPerson. "Seidl" 16000)))
```

```
Dr. Petter
HiWi Stefan
Prof. Seidl
```

# Multidispatching

## Pro

- Generalization of an established technique
- Directly solves problem
- Eliminates boilerplate code
- Compatible with modular compilation/type checking

## Con

- Counters privileged 1st parameter
- Runtime overhead
- New exceptions when used with multi-inheritance
- *Most Specific Method* ambiguous

## Other Solutions (extract)

- Dylan
- Scala

## Lessons Learned

1. Dynamically dispatched methods are complex interaction of static and dynamic techniques
2. Single Dispatching as in major OO-Languages
3. Making use of Open Source Compilers
4. Multi Dispatching generalizes single dispatching
5. Multi Dispatching Perl6
6. Multi Dispatching Clojure

Section 6

**Further materials**

# Further reading...

[1] hotspot/src/share/vm/interpreter/linkResolver.cpp.
OpenJDK 7 Hotspot JIT VM.
http://hg.openjdk.java.net/jdk7/jdk7.

[2] jdk/src/share/classes/sun/tools/java/ClassDefinition.java.
OpenJDK 7 Javac.
http://hg.openjdk.java.net/jdk7/jdk7.

[3] C. Clifton, T. Millstein, G. T. Leavens, and C. Chambers.
Multijava: Design rationale, compiler implementation, and applications.
*ACM Transactions on Programming Languages and Systems (TOPLAS)*, May 2006.

[4] J. Gosling, B. Joy, G. Steele, and G. Bracha.
*The Java Language Specification, Third Edition.*
Addison-Wesley Longman, Amsterdam, 3 edition, June 2005.

[5] S. Halloway.
*Programming Clojure.*
Pragmatic Bookshelf, 1st edition, 2009.

[6] T. Lindholm, F. Yellin, G. Bracha, and A. Buckley.
*The Java Virtual Machine Specification.*
Addison-Wesley Professional, Java SE7 edition, 2013.

[7] R. Muschevici, A. Potanin, E. Tempero, and J. Noble.
Multiple dispatch in practice.
*23rd ACM SIGPLAN conference on Object-oriented programming systems languages and applications (OOPSLA)*, September 2008.

# **Programming Languages**

Multiple Inheritance

Dr. Michael Petter
Winter term 2019

# Outline

# Outline

"Wouldn't it be nice to inherit from several parents?"

# Interface vs. Implementation inheritance

The classic motivation for inheritance is implementation inheritance

- *Code reusage*
- Child specializes parents, replacing particular methods with custom ones
- Parent acts as library of common behaviours
- Implemented in languages like C++ or Lisp

Code sharing in interface inheritance inverts this relation

- *Behaviour contract*
- Child provides methods, with signatures predetermined by the parent
- Parent acts as generic code frame with room for customization
- Implemented in languages like Java or C#

# Interface Inheritance

"So how do we lay out objects in memory anyway?"

## Excursion: Brief introduction to LLVM IR

LLVM intermediate representation as reference semantics:

```
;(recursive) struct definitions
%struct.A = type { i32, %struct.B, i32(i32)* }
%struct.B = type { i64, [10 x [20 x i32]], i8 }

;(stack-) allocation of objects
%a = alloca %struct.A
;adress computation for selection in structure (pointers):
%1 = getelementptr %struct.A* %a, i64 0, i64 2
;load from memory
%2 = load i32(i32)* %1
;indirect call
%retval = call i32 (i32)* %2(i32 42)
```

Retrieve the memory layout of a compilation unit with:
```
clang -cc1 -x c++ -v -fdump-record-layouts -emit-llvm source.cpp
```
Retrieve the IR Code of a compilation unit with:
```
clang -O1 -S -emit-llvm source.cpp -o IR.llvm
```

```
class A {
  int a; int f(int);
};
class B : public A {
  int b; int g(int);
};
class C : public B {
  int c; int h(int);
};

...

C c;
c.g(42);
```



C (=A/B)

int a
int b
int c

```
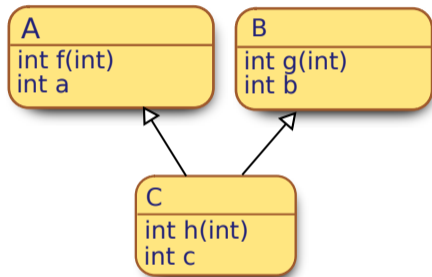%class.C = type { %class.B, i32 }
%class.B = type { %class.A, i32 }
%class.A = type { i32 }
```

```
%c = alloca %class.C
%1 = bitcast %class.C* %c to %class.B*
%2 = call i32 @_g(%class.B* %1, i32 42) ; g is statically known
```

```
class A {
  int a; int f(int);
};
class B : public A {
  int b; int g(int);
};
class C : public B {
  int c; int h(int);
};
int B::g(int p) {
  return p+b;
};
```

**C**

int a
int b
int c

```
%class.C = type { %class.B, i32 }
%class.B = type { %class.A, i32 }
%class.A = type { i32 }
```

```
define i32 @_g(%class.B* %this, i32 %p) {
  %1 = getelementptr %class.B* %this, i64 0, i32 1
  %2 = load i32* %1
  %3 = add i32 %2, %p
  ret i32 %3
}
```

"Now what about polymorphic calls?"

# Single-Dispatching implementation choices

Single-Dispatching needs runtime action:

1. Manual search run through the super-chain (Java Interpreter ⤳ last talk)

```
call i32 @__dispatch(%class.C* %c,i32 42,i32* "f(int,void)")
```

# Single-Dispatching implementation choices

Single-Dispatching needs runtime action:

1. Manual search run through the super-chain (Java Interpreter $\rightsquigarrow$ last talk)

```
call i32 @__dispatch(%class.C* %c,i32 42,i32* "f(int,void)")
```

2. Caching the dispatch result ($\rightsquigarrow$ Hotspot/JIT)

```
; caching the recent result value of the __dispatch function
; call i32 @__dispatch(%class.C* %c,i32 42)
assert (%c type %class.D) ; verify objects class presumption
call i32 @_f_from_D(%class.C* %c, i32 42) ; directly call f
```

# Single-Dispatching implementation choices

Single-Dispatching needs runtime action:

1. Manual search run through the super-chain (Java Interpreter ⤳ last talk)

```
call i32 @__dispatch(%class.C* %c,i32 42,i32* "f(int,void)")
```

2. Caching the dispatch result (⤳ Hotspot/JIT)

```
; caching the recent result value of the __dispatch function
; call i32 @__dispatch(%class.C* %c,i32 42)
assert (%c type %class.D) ; verify objects class presumption
call i32 @_f_from_D(%class.C* %c, i32 42) ; directly call f
```

3. Precomputing the dispatching result in tables

# Single-Dispatching implementation choices

Single-Dispatching needs runtime action:

1. Manual search run through the super-chain (Java Interpreter $\rightsquigarrow$ last talk)

```
call i32 @__dispatch(%class.C* %c,i32 42,i32* "f(int,void)")
```

2. Caching the dispatch result ($\rightsquigarrow$ Hotspot/JIT)

```
; caching the recent result value of the __dispatch function
; call i32 @__dispatch(%class.C* %c,i32 42)
assert (%c type %class.D) ; verify objects class presumption
call i32 @_f_from_D(%class.C* %c, i32 42) ; directly call f
```

3. Precomputing the dispatching result in tables
   1. Full 2-dim matrix

|   | f() | g() | h() | i() | j() | k() | l() | m() | n() |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| A | 1   |     |     |     |     |     |     |     |     |
| B | 1   | 2   |     |     |     |     |     |     |     |
| C | 3   |     | 4   |     |     |     |     |     |     |
| D | 3   | 2   | 4   | 5   |     |     |     |     |     |
| E |     |     |     |     |     | 6   |     | 7   |     |
| F |     |     |     |     | 8   | 9   |     | 7   |     |

# Single-Dispatching implementation choices

Single-Dispatching needs runtime action:

1. Manual search run through the super-chain (Java Interpreter ↝ last talk)

```
call i32 @__dispatch(%class.C* %c,i32 42,i32* "f(int,void)")
```

2. Caching the dispatch result (↝ Hotspot/JIT)

```
; caching the recent result value of the __dispatch function
; call i32 @__dispatch(%class.C* %c,i32 42)
assert (%c type %class.D) ; verify objects class presumption
call i32 @_f_from_D(%class.C* %c, i32 42) ; directly call f
```

3. Precomputing the dispatching result in tables
   1. Full 2-dim matrix
   2. 1-dim Row Displacement Dispatch Tables

|   | f() | g() | h() | i() | j() | k() | l() | m() | n() |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| A | 1   |     |     |     |     |     |     |     |     |
| B | 1   | 2   |     |     |     |     |     |     |     |
| C | 3   |     | 4   |     |     |     |     |     |     |
| D | 3   | 2   | 4   | 5   |     |     |     |     |     |
| E |     |     |     |     |     | 6   |     | 7   |     |
| F |     |     |     |     | 8   | 9   |     | 7   |     |

|   | A | B |   |   | F |   |   |   |
|---|---|---|---|---|---|---|---|---|
|   | 1 | 1 | 2 | ... | 8 | 9 |   | 7 |

# Single-Dispatching implementation choices

Single-Dispatching needs runtime action:

**1** Manual search run through the super-chain (Java Interpreter ⤳ last talk)

```
call i32 @__dispatch(%class.C* %c,i32 42,i32* "f(int,void)")
```

**2** Caching the dispatch result (⤳ Hotspot/JIT)

```
; caching the recent result value of the __dispatch function
; call i32 @__dispatch(%class.C* %c,i32 42)
assert (%c type %class.D) ; verify objects class presumption
call i32 @_f_from_D(%class.C* %c, i32 42) ; directly call f
```

**3** Precomputing the dispatching result in tables
   **1** Full 2-dim matrix
   **2** 1-dim Row Displacement Dispatch Tables
   **3** Virtual Tables
      (⤳ LLVM/GNU C++,this talk)

|   | f() | g() | h() | i() | j() | k() | l() | m() | n() |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| A | 1   |     |     |     |     |     |     |     |     |
| B | 1   | 2   |     |     |     |     |     |     |     |
| C | 3   |     | 4   |     |     |     |     |     |     |
| D | 3   | 2   | 4   | 5   |     |     |     |     |     |
| E |     |     |     |     |     | 6   |     | 7   |     |
| F |     |     |     |     | 8   | 9   |     | 7   |     |

| A | B |   |   |   |   | F |   |   |   |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 2 | ... |   | 8 | 9 |   | 7 |   |

# Object layout – virtual methods

```cpp
class A {
  int a; virtual int f(int);
         virtual int g(int);
         virtual int h(int);
};
class B : public A {
  int b; int g(int);
};
class C : public B {
  int c; int h(int);
}; ...
C c;
c.g(42);
```



```
%class.C = type { %class.B, i32, [4 x i8] }
%class.B = type { [12 x i8], i32 }
%class.A = type { i32 (...)**, i32 }
```

```llvm
%c.vptr = bitcast %class.C* %c to i32 (%class.B*, i32)*** ; vtbl
%1 = load (%class.B*, i32)*** %c.vptr      ; dereference vptr
%2 = getelementptr %1, i64 1               ; select g()-entry
%3 = load (%class.B*, i32)** %2            ; dereference g()-entry
%4 = call i32 %3(%class.B* %c, i32 42)
```

"So how do we include several parent objects?"

A

int f(int)
int a

B

int g(int)
int b

C

int h(int)
int c

# Static Type Casts

```cpp
class A {
  int a; int f(int);
};
class B {
  int b; int g(int);
};
class C : public A , public B {
  int c; int h(int);
};
...
B* b = new C();
```



```llvm
%class.C = type { %class.A, %class.B, i32 }
%class.A = type { i32 }
%class.B = type { i32 }
```

```llvm
%1 = call i8* @_new(i64 12)
call void @_memset.p0i8.i64(i8* %1, i8 0, i64 12, i32 4, i1 false)
%2 = getelementptr i8* %1, i64 4        ; select B-offset in C
%b = bitcast i8* %2 to %class.B*
```

# Static Type Casts

```
class A {
  int a; int f(int);
};
class B {
  int b; int g(int);
};
class C : public A , public B {
  int c; int h(int);
};
...
B* b = new C();
```



```
%class.C = type { %class.A, %class.B, i32 }
%class.A = type { i32 }
%class.B = type { i32 }
```

```
%1 = call i8* @_new(i64 12)
call void @_memset.p0i8.i64(i8* %1, i8 0, i64 12, i32 4, i1 false)
%2 = getelementptr i8* %1, i64 4          ; select B-offset in C
%b = bitcast i8* %2 to %class.B*
```

⚠ implicit casts potentially add a constant to the object pointer.

# Static Type Casts

```cpp
class A {
  int a; int f(int);
};
class B {
  int b; int g(int);
};
class C : public A , public B {
  int c; int h(int);
};
...
B* b = new C();
```



```llvm
%class.C = type { %class.A, %class.B, i32 }
%class.A = type { i32 }
%class.B = type { i32 }
```

```llvm
%1 = call i8* @_new(i64 12)
call void @_memset.p0i8.i64(i8* %1, i8 0, i64 12, i32 4, i1 false)
%2 = getelementptr i8* %1, i64 4          ; select B-offset in C
%b = bitcast i8* %2 to %class.B*
```

⚠ implicit casts potentially add a constant to the object pointer.

⚠ getelementptr implements $\triangle B$ as $4 \cdot i8$!

# Keeping Calling Conventions

```
class A {
    int a; int f(int);
};
class B {
    int b; int g(int);
};
class C : public A , public B {
    int c; int h(int);
};
...
C c;
c.g(42);
```



```
%class.C = type { %class.A, %class.B, i32 }
%class.A = type { i32 }
%class.B = type { i32 }
```

```
%c = alloca %class.C
%1 = bitcast %class.C* %c to i8*
%2 = getelementptr i8* %1, i64 4       ; select B-offset in C
%3 = call i32 @_g(%class.B* %2, i32 42) ; g is statically known
```

```cpp
class A { void f(int); };
class B { void f(int); };
class C : public A, public B {};

C* pc;
pc->f(42);
```

⚠ Which method is called?

Solution I: Explicit qualification

```cpp
pc->A::f(42);
pc->B::f(42);
```

Solution II: Automagical resolution

Idea: The Compiler introduces a linear order on the nodes of the inheritance graph

# Linearization

| **Principle 1: Inheritance Relation** | **Principle 2: Multiplicity Relation** |
|---|---|
| Defined by parent-child. Example: $C(A, B) \implies C \mathrel{-\!\!\triangleright} A \wedge C \mathrel{-\!\!\triangleright} B$ $\qquad\qquad \longrightarrow\!\!\triangleright$ | Defined by the succession of multiple parents. Example: $C(A, B) \implies A \rightarrow B$ |

In General:

1. Inheritance is a uniform mechanism, and its searches ($\rightarrow$ total order) apply identically for all object fields or methods

2. In the literature, we also find the set of constraints to create a linearization as <u>M</u>ethod <u>R</u>esolution <u>O</u>rder

3. Linearization is a best-effort approach at best

# MRO via DFS

$A(B, C)\ B(W)\ C(W)$

# MRO via DFS

## Leftmost Preorder Depth-First Search

$$L[A] = ABWC$$

⚠ Principle 1 *inheritance* is violated

  Python: classical python objects ($\leq 2.1$) use LPDFS!

## LPDFS with Duplicate Cancellation



$A(B, C)\ B(W)\ C(W)$

# MRO via DFS



$A(B, C)$ $B(W)$ $C(W)$



$A(B, C)B(V, W)C(W, V)$

## Leftmost Preorder Depth-First Search

$$L[A] = ABWC$$

⚠ Principle 1 *inheritance* is violated

Python: classical python objects ($\leq 2.1$) use LPDFS!

## LPDFS with Duplicate Cancellation

$$L[A] = ABCW$$

✓ Principle 1 *inheritance* is fixed

Python: new python objects (2.2) use LPDFS(DC)!

## LPDFS with Duplicate Cancellation

## MRO via DFS



$A(B, C) \ B(W) \ C(W)$



$A(B, C)B(V, W)C(W, V)$

### Leftmost Preorder Depth-First Search

$$L[A] = ABWC$$

⚠ Principle 1 *inheritance* is violated

Python: classical python objects ($\leq 2.1$) use LPDFS!

### LPDFS with Duplicate Cancellation

$$L[A] = ABCW$$

✓ Principle 1 *inheritance* is fixed

Python: new python objects (2.2) use LPDFS(DC)!

### LPDFS with Duplicate Cancellation

$$L[A] = ABCWV$$

⚠ Principle 2 *multiplicity* not fulfillable
⚠ However $B \rightarrow C \implies W \rightarrow V$??

# MRO via Refined Postorder DFS

$A(B, C) \ B(F, D) \ C(E, H)$
$D(G) \ E(G) \ F(W) \ G(W) \ H(W)$

# MRO via Refined Postorder DFS



## Reverse Postorder Rightmost DFS

$$L[A] = ABFDCEGHW$$

✓ Linear extension of inheritance relation

## RPRDFS



$A(B, C)\ B(F, D)\ C(E, H)$
$D(G)\ E(G)\ F(W)\ G(W)\ H(W)$



$A(B, C)\ B(F, G)\ C(D, E)$

# MRO via Refined Postorder DFS



**Reverse Postorder Rightmost DFS**

$$L[A] = ABFDCEGHW$$

✓ Linear extension of inheritance relation

$A(B, C) \; B(F, D) \; C(E, H)$
$D(G) \; E(G) \; F(W) \; G(W) \; H(W)$

**RPRDFS**

$$L[A] = ABCDGEF$$

⚠ But principle 2 *multiplicity* is violated!



$A(B, C) \; B(F, G) \; C(D, E)$

# MRO via Refined Postorder DFS



$A(B, C)\; B(F, D)\; C(E, H)$
$D(G)\; E(G)\; F(W)\; G(W)\; H(W)$

## Reverse Postorder Rightmost DFS

$$L[A] = ABFDCEGHW$$

✓ Linear extension of inheritance relation

## RPRDFS

$$L[A] = ABCDGEF$$

⚠ But principle 2 *multiplicity* is violated!

## Refined RPRDFS



$A(B, C)\; B(F, G)\; C(D, E)$

# MRO via Refined Postorder DFS



## Reverse Postorder Rightmost DFS

$$L[A] = ABFDCEGHW$$

✓ Linear extension of inheritance relation

$A(B, C)\ B(F, D)\ C(E, H)$
$D(G)\ E(G)\ F(W)\ G(W)\ H(W)$

## RPRDFS

$$L[A] = ABCDGEF$$

⚠ But principle 2 *multiplicity* is violated!

CLOS: uses Refined RPDFS [3]

## Refined RPRDFS

$$L[A] = ABCDEFG$$

✓ Refine graph with conflict edge & rerun RPRDFS!



$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

**Extension Principle: Monotonicity**

If $C_1 \to C_2$ in $C$'s linearization, then $C_1 \to C_2$ for every linearization of $C$'s children.

$A(B, C) \; B(F, G) \; C(D, E)$
$D(G) \; E(F)$

## Refined RPRDFS

⚠ *Monotonicity* is not guaranteed!

## Extension Principle: Monotonicity

If $C_1 \to C_2$ in $C$'s linearization, then $C_1 \to C_2$ for every linearization of $C$'s children.



$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

# MRO via Refined Postorder DFS

**Refined RPRDFS**

⚠ *Monotonicity* is not guaranteed!

**Extension Principle: Monotonicity**

If $C_1 \to C_2$ in $C$'s linearization, then $C_1 \to C_2$ for every linearization of $C$'s children.

$$L[A] = A\ B\ C\ D\ E\ F\ G \quad \implies \quad F \to G$$

$$L[C] = C\ D\ G\ E\ F \quad \implies \quad G \to F$$



$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

A linearization $L$ is an attribute $L[C]$ of a class $C$. Classes $B_1, \ldots, B_n$ are superclasses to child class $C$, defined in the *local precedence order* $C(B_1 \ldots B_n)$. Then

$$L[C] = C \cdot \bigsqcup (L[B_1], \ldots, L[B_n], B_1 \cdot \ldots \cdot B_n) \quad | \quad C(B_1, \ldots, B_n)$$

$$L[Object] = Object$$

with

$$\bigsqcup_i (L_i) = \begin{cases} c \cdot (\bigsqcup_i (L_i \setminus c)) & \text{if } \exists_{\min k} \forall_j \ c = head(L_k) \notin tail(L_j) \\ \triangle \text{ fail} & \text{else} \end{cases}$$

$$A(B, C)\ B(F, G)\ C(D, E)$$
$$D(G)\ E(F)$$

$$
\begin{array}{ll}
L[G] & G \\
L[F] & F \\
L[E] & \\
L[D] & \\
L[B] & \\
L[C] & \\
L[A] & \\
\end{array}
$$

$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

$$
\begin{array}{ll}
L[G] & G \\
L[F] & F \\
L[E] & E \cdot F \\
L[D] & D \cdot G \\
L[B] & \\
L[C] & \\
L[A] &
\end{array}
$$

$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

$L[G] \quad G$
$L[F] \quad F$
$L[E] \quad E \cdot F$
$L[D] \quad D \cdot G$
$L[B] \quad B \cdot (L[F] \sqcup L[G] \sqcup (F \cdot G))$
$L[C]$
$L[A]$

$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

$L[G]$  $G$
$L[F]$  $F$
$L[E]$  $E \cdot F$
$L[D]$  $D \cdot G$
$L[B]$  $B \cdot (F \sqcup G \sqcup (F \cdot G))$
$L[C]$
$L[A]$

$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

$$
\begin{array}{ll}
L[G] & G \\
L[F] & F \\
L[E] & E \cdot F \\
L[D] & D \cdot G \\
L[B] & B \cdot F \cdot G \\
L[C] & \\
L[A] &
\end{array}
$$

$L[G]$    $G$
$L[F]$    $F$
$L[E]$    $E \cdot F$
$L[D]$    $D \cdot G$
$L[B]$    $B \cdot F \cdot G$
$L[C]$    $C \cdot (L[D] \sqcup L[E] \sqcup (D \cdot E))$
$L[A]$

$A(B, C) \ B(F, G) \ C(D, E)$
$D(G) \ E(F)$

$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

$L[G] \quad G$
$L[F] \quad F$
$L[E] \quad E \cdot F$
$L[D] \quad D \cdot G$
$L[B] \quad B \cdot F \cdot G$
$L[C] \quad C \cdot ((D \cdot G) \sqcup (E \cdot F) \sqcup (D \cdot E))$
$L[A]$

$A(B, C) \; B(F, G) \; C(D, E)$
$D(G) \; E(F)$

$$
\begin{array}{ll}
L[G] & G \\
L[F] & F \\
L[E] & E \cdot F \\
L[D] & D \cdot G \\
L[B] & B \cdot F \cdot G \\
L[C] & C \cdot D \cdot (G \sqcup (E \cdot F) \sqcup E) \\
L[A] &
\end{array}
$$

$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

$L[G]$    $G$
$L[F]$    $F$
$L[E]$    $E \cdot F$
$L[D]$    $D \cdot G$
$L[B]$    $B \cdot F \cdot G$
$L[C]$    $C \cdot D \cdot G \cdot E \cdot F$
$L[A]$

$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

$L[G] \quad G$

$L[F] \quad F$

$L[E] \quad E \cdot F$

$L[D] \quad D \cdot G$

$L[B] \quad B \cdot F \cdot G$

$L[C] \quad C \cdot D \cdot G \cdot E \cdot F$

$L[A] \quad A \cdot ((B \cdot F \cdot G) \sqcup (C \cdot D \cdot G \cdot E \cdot F) \sqcup (B \cdot C))$

$A(B, C)\ B(F, G)\ C(D, E)$
$D(G)\ E(F)$

$L[G]$  $G$

$L[F]$  $F$

$L[E]$  $E \cdot F$

$L[D]$  $D \cdot G$

$L[B]$  $B \cdot F \cdot G$

$L[C]$  $C \cdot D \cdot G \cdot E \cdot F$

$L[A]$  $A \cdot B \cdot C \cdot D \cdot ((F \cdot G) \sqcup (G \cdot E \cdot F))$

$$A(B, C)\ B(F, G)\ C(D, E)$$
$$D(G)\ E(F)$$

| | |
|---|---|
| $L[G]$ | $G$ |
| $L[F]$ | $F$ |
| $L[E]$ | $E \cdot F$ |
| $L[D]$ | $D \cdot G$ |
| $L[B]$ | $B \cdot F \cdot G$ |
| $L[C]$ | $C \cdot D \cdot G \cdot E \cdot F$ |
| $L[A]$ | ⚠ fail |

$L[G]$   $G$
$L[F]$   $F$
$L[E]$   $E \cdot F$
$L[D]$   $D \cdot G$
$L[B]$   $B \cdot F \cdot G$
$L[C]$   $C \cdot D \cdot G \cdot E \cdot F$
$L[A]$   ⚠ fail



$A(B, C) \; B(F, G) \; C(D, E)$
$D(G) \; E(F)$

C3 detects and reports a violation of *monotonicity* with the addition of A(B,C) to the class set.
C3 linearization [1]: is used in *Python 3*, *Perl 6*, and *Solidity*

# Linearization vs. explicit qualification

## Linearization

- No switch/duplexer code necessary
- No explicit naming of qualifiers
- Unique `super` reference
- Reduces number of multi-dispatching conflicts

## Qualification

- More flexible, fine-grained
- Linearization choices may be awkward or unexpected

## Languages with automatic linearization exist

- *CLOS* Common Lisp Object System
- *Solidity, Python 3* and *Perl 6* with C3
- Prerequisite for → Mixins

"And what about dynamic dispatching in Multiple Inheritance?"

# Virtual Tables for Multiple Inheritance

```cpp
class A {
  int a; virtual int f(int);
};
class B {
  int b; virtual int f(int);
          virtual int g(int);
};
class C : public A , public B {
  int c; int f(int);
};
...
C c;
B* pb = &c;
pb->f(42);
```



```llvm
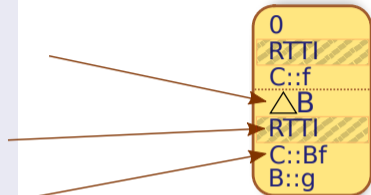%class.C = type { %class.A, [12 x i8], i32 }
%class.A = type { i32 (...)**, i32 }
%class.B = type { i32 (...)**, i32 }
```

```llvm
; B* pb = &c;
%0 = bitcast %class.C* %c to i8*      ; type fumbling
%1 = getelementptr i8* %0, i64 16     ; offset of B in C
%2 = bitcast i8* %1 to %class.B*      ; get typing right
store %class.B* %2, %class.B** %pb    ; store to pb
```

# Virtual Tables for Multiple Inheritance

```cpp
class A {
  int a; virtual int f(int);
};
class B {
  int b; virtual int f(int);
          virtual int g(int);
};
class C : public A , public B {
  int c; int f(int);
};
...
C c;
B* pb = &c;
pb->f(42);
```



```
%class.C = type { %class.A, [12 x i8], i32 }
%class.A = type { i32 (...)**, i32 }
%class.B = type { i32 (...)**, i32 }
```

```
; pb->f(42);
%0 = load %class.B** %pb                              ;load the b-pointer
%1 = bitcast %class.B* %0 to i32 (%class.B*, i32)***  ;cast to vtable
%2 = load i32(%class.B*, i32)*** %1                   ;load vptr
%3 = getelementptr i32 (%class.B*, i32)** %2, i64 0   ;select f() entry
%4 = load i32(%class.B*, i32)** %3                    ;load function pointer
%5 = call i32 %4(%class.B* %0, i32 42)
```

# Virtual Tables for Multiple Inheritance

```
class A {
  int a; virtual int f(int);
};
class B {
  int b; virtual int f(int);
         virtual int g(int);
};
class C : public A , public B {
  int c; int f(int);
};
...
C c;
B* pb = &c;
pb->f(42);
```



```
%class.C = type { %class.A, [12 x i8], i32 }
%class.A = type { i32 (...)**, i32 }
%class.B = type { i32 (...)**, i32 }
```

```
; pb->f(42);
%0 = load %class.B** %pb                              ;load the b-pointer
%1 = bitcast %class.B* %0 to i32 (%class.B*, i32)***  ;cast to vtable
%2 = load i32(%class.B*, i32)*** %1                   ;load vptr
%3 = getelementptr i32 (%class.B*, i32)** %2, i64 0   ;select f() entry
%4 = load i32(%class.B*, i32)** %3                    ;load function pointer
%5 = call i32 %4(%class.B* %0, i32 42)
```

**A Basic Virtual Table**

consists of different parts:

1. *offset to top* of an enclosing objects memory representation

2. *typeinfo pointer* to an RTTI object
   (not relevant for us)

3. *virtual function pointers* for resolving virtual methods

```
0
RTTI
C::f
△B
RTTI
C::Bf
B::g
```

- Virtual tables are composed when multiple inheritance is used
- The `vptr` fields in objects are pointers to their corresponding virtual-subtables
- Casting preserves the link between an object and its corresponding virtual-subtable
- `clang -cc1 -fdump-vtable-layouts -emit-llvm code.cpp` yields the vtables of a compilation unit

```
class A { int a; };
class B { virtual int f(int); };
class C : public A , public B {
  int c; int f(int);
};
C* c = new C();
c->f(42);
```

```
B* b = new C();
b->f(42);
```

# Casting Issues

```
class A { int a; };
class B { virtual int f(int); };
class C : public A , public B {
  int c; int f(int);
};
C* c = new C();
c->f(42);
```

```
B* b = new C();
b->f(42);
```



⚠ `this`-Pointer for `C::f` is expected to point to `C`

# Casting Issues

```
class A { int a; };
class B { virtual int f(int); };
class C : public A , public B {
  int c; int f(int);
};
C* c = new C();
c->f(42);
```

```
B* b = new C();
b->f(42);
```



```
0
RTTI
C::f
△B
RTTI
C::Bf
```

⚠ `this`-Pointer for `C::f` is
expected to point to `C`

C::Bf

C::f

# Thunks

**Solution: *thunk*s**

. . . are trampoline methods, delegating the virtual method to its original implementation with an adapted `this`-reference

```
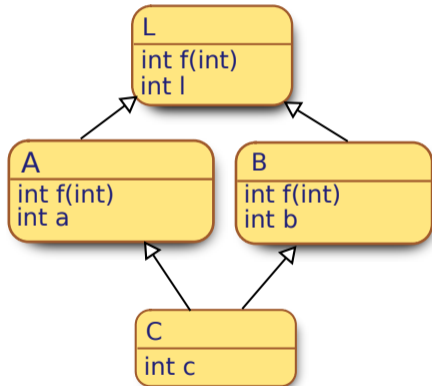define i32 @__f(%class.B* %this, i32 %i) {
  %1 = bitcast %class.B* %this to i8*
  %2 = getelementptr i8* %1, i64 -16      ; sizeof(A)=16
  %3 = bitcast i8* %2 to %class.C*
  %4 = call i32 @_f(%class.C* %3, i32 %i)
  ret i32 %4
}
```

↝ `B-in-C`-vtable entry for `f(int)` is the thunk `_f(int)`

# Thunks

**Solution:** *thunk*s

. . . are trampoline methods, delegating the virtual method to its original implementation with an adapted `this`-reference

```
define i32 @__f(%class.B* %this, i32 %i) {
  %1 = bitcast %class.B* %this to i8*
  %2 = getelementptr i8* %1, i64 -16      ; sizeof(A)=16
  %3 = bitcast i8* %2 to %class.C*
  %4 = call i32 @_f(%class.C* %3, i32 %i)
  ret i32 %4
}
```

⤳ `B-in-C`-vtable entry for `f(int)` is the thunk `_f(int)`

⤳ `_f(int)` adds a compiletime constant $\Delta B$ to `this` before calling `f(int)`

# Thunks

## Solution: *thunk*s

. . . are trampoline methods, delegating the virtual method to its original implementation with an adapted `this`-reference

```
define i32 @__f(%class.B* %this, i32 %i) {
  %1 = bitcast %class.B* %this to i8*
  %2 = getelementptr i8* %1, i64 -16      ; sizeof(A)=16
  %3 = bitcast i8* %2 to %class.C*
  %4 = call i32 @_f(%class.C* %3, i32 %i)
  ret i32 %4
}
```

$\rightsquigarrow$ B-in-C-vtable entry for f(int) is the thunk _f(int)

$\rightsquigarrow$ _f(int) adds a compiletime constant $\Delta B$ to `this` before calling f(int)

$\rightsquigarrow$ f(int) addresses its locals relative to what it assumes to be a C pointer

"But what if there are common ancestors?"

Standard C++ multiple inheritance conceptually duplicates representations for common ancestors:

# Common Bases – Duplicated Bases

Standard C++ multiple inheritance conceptually duplicates representations for common ancestors:

# Duplicated Base Classes

```cpp
class L {
  int l; virtual void f(int);
};
class A : public L {
  int a; void f(int);
};
class B : public L {
  int b; void f(int);
};
class C : public A , public B {
  int c;
};
...
C c;
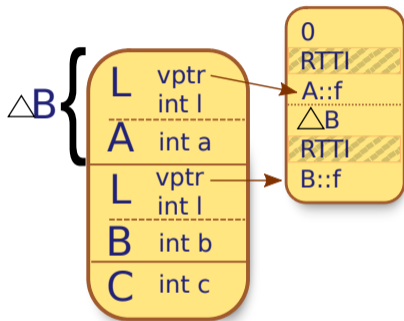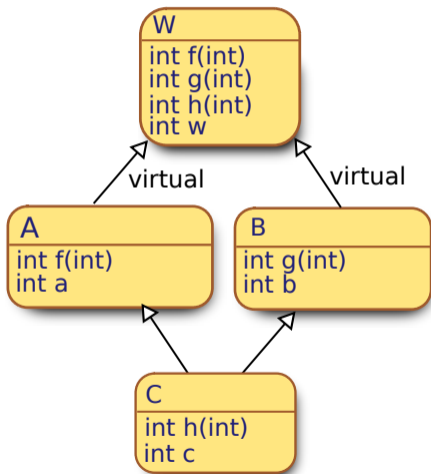L* pl = &c;
pl->f(42); // where to dispatch?
C* pc = (C*)pl;
```



```llvm
%class.C = type { %class.A, %class.B,
                  i32, [4 x i8] }
%class.A = type { [12 x i8], i32 }
%class.B = type { [12 x i8], i32 }
%class.L = type { i32 (...)**, i32 }
```

⚠ Ambiguity!

# Duplicated Base Classes

```cpp
class L {
  int l; virtual void f(int);
};
class A : public L {
  int a; void f(int);
};
class B : public L {
  int b; void f(int);
};
class C : public A , public B {
  int c;
};
...
C c;
L* pl = (B*)&c;
pl->f(42); // where to dispatch?
C* pc = (C*)pl;
```



```llvm
%class.C = type { %class.A, %class.B,
                  i32, [4 x i8] }
%class.A = type { [12 x i8], i32 }
%class.B = type { [12 x i8], i32 }
%class.L = type { i32 (...)**, i32 }
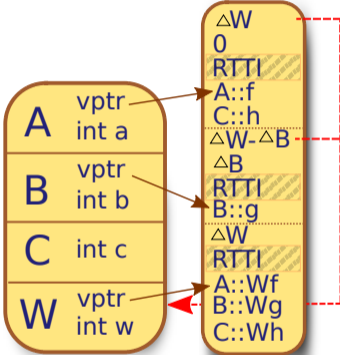```

⚠ Ambiguity!

# Duplicated Base Classes

```
class L {
  int l; virtual void f(int);
};
class A : public L {
  int a; void f(int);
};
class B : public L {
  int b; void f(int);
};
class C : public A , public B {
  int c;
};
...
C c;
L* pl = (B*)&c;
pl->f(42); // where to dispatch?
C* pc = (C*)(B*)pl;
```



```
%class.C = type { %class.A, %class.B,
                  i32, [4 x i8] }
%class.A = type { [12 x i8], i32 }
%class.B = type { [12 x i8], i32 }
%class.L = type { i32 (...)**, i32 }
```

# Common Bases – Shared Base Class

Optionally, C++ multiple inheritance enables a shared representation for common ancestors, creating the *diamond pattern*:

# Shared Base Class

```cpp
class W {
  int w; virtual void f(int);
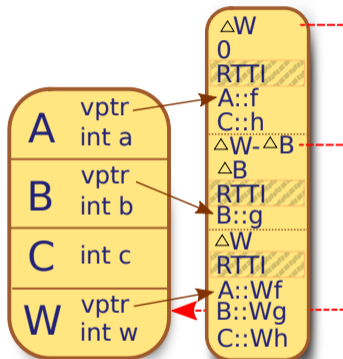  virtual void g(int);
  virtual void h(int);
};
class A : public virtual W {
  int a; void f(int);
};
class B : public virtual W {
  int b; void g(int);
};
class C : public A, public B {
  int c; void h(int);
};
...
C* pc;
pc->f(42);
```



⚠ Ambiguities
⤳ e.g. overriding f in A *and* B

# Shared Base Class

```cpp
class W {
  int w; virtual void f(int);
  virtual void g(int);
  virtual void h(int);
};
class A : public virtual W {
  int a; void f(int);
};
class B : public virtual W {
  int b; void g(int);
};
class C : public A, public B {
  int c; void h(int);
};
...
C* pc;
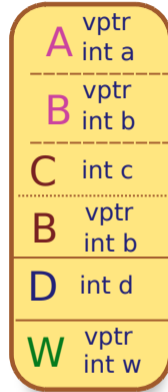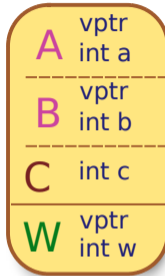pc->B::f(42);
((W*)pc)->h(42);
((B*)pc)->f(42);
```



⚠ Offsets to virtual base

# Dynamic Type Casts

```
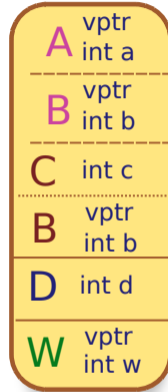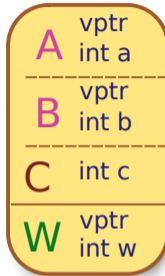class A : public virtual W {
...
};
class B : public virtual W {
...
};
class C : public A , public B {
...
};
class D : public C,
          public B {
...
};
...
C c;
W* pw = &c;
C* pc = (C*)pw; // Compile error
```



⚠ No guaranteed *constant* offsets between virtual bases and subclasses ⤳ No static casting!

# Dynamic Type Casts

```
class A : public virtual W {
...
};
class B : public virtual W {
...
};
class C : public A , public B {
...
};
class D : public C,
          public B {
...
};
...
C c;
W* pw = &c;
C* pc = dynamic_cast<C*>(pw);
```



⚠ No guaranteed *constant* offsets between virtual bases and subclasses ⇝ No static casting!
⚠ *Dynamic casting* makes use of *offset-to-top*

```
class W { virtual int f(int); };
class A : virtual W { int a; };
class B : virtual W { int b; };
class C : public A , public B {
  int c; int f(int);
};
B* b = new C();
b->f(42);
```

```
W* w = new C();
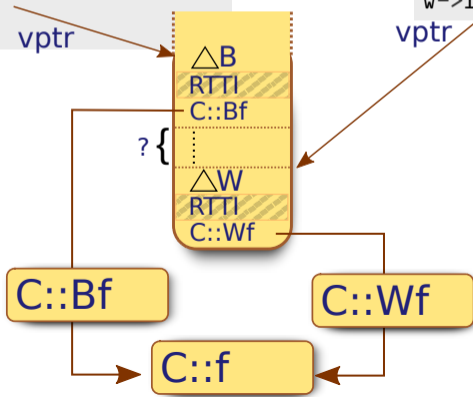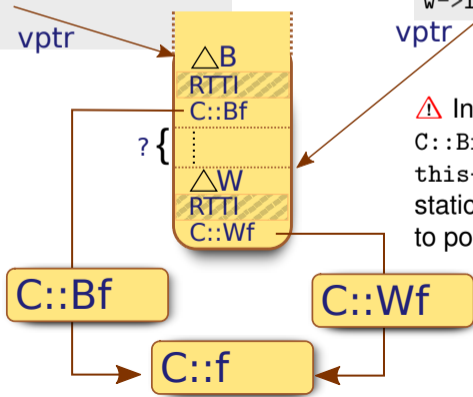w->f(42);
```

# Again: Casting Issues

```
class W { virtual int f(int); };
class A : virtual W { int a; };
class B : virtual W { int b; };
class C : public A , public B {
  int c; int f(int);
};
B* b = new C();
b->f(42);
```

```
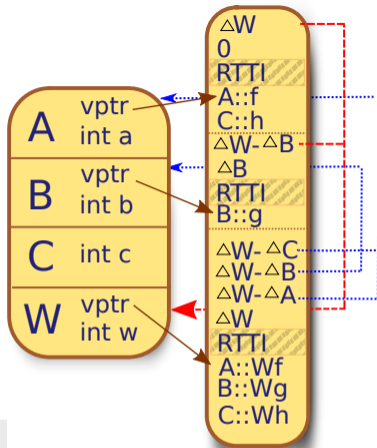W* w = new C();
w->f(42);
```

vptr

vptr



⚠ In a conventional thunk `C::Bf` adjusts the **this**-pointer with a statically known constant to point to `C`

# Virtual Thunks

```cpp
class W { ...
virtual void g(int);
};
class A : public virtual W {...};
class B : public virtual W {
  int b; void g(int i){ };
};
class C : public A,public B{...};
C c;
W* pw = &c;
pw->g(42);
```



```llvm
define void @__g(%class.B* %this, i32 %i) { ; virtual thunk to B::g
  %1 = bitcast %class.B* %this to i8*
  %2 = bitcast i8* %1 to i8**
  %3 = load i8** %2              ; load W-vtable ptr
  %4 = getelementptr i8* %3, i64 -32 ; -32 bytes is g-entry in vcalls
  %5 = bitcast i8* %4 to i64*
  %6 = load i64* %5             ; load g's vcall offset
  %7 = getelementptr i8* %1, i64 %6  ; navigate to vcalloffset+ Wtop
  %8 = bitcast i8* %7 to %class.B*
  call void @_g(%class.B* %8, i32 %i)
  ret void
}
```

# Virtual Tables for Virtual Bases (⇝ C++-ABI)

**A Virtual Table for a Virtual Subclass**

gets a *virtual base pointer*

**A Virtual Table for a Virtual Base**

consists of different parts:

1. *virtual call offsets* per virtual function for adjusting `this` dynamically
2. *offset to top* of an enclosing objects heap representation
3. *typeinfo pointer* to an RTTI object (not relevant for us)
4. *virtual function pointers* for resolving virtual methods



△W-△B

△B

RTTI
B::g

△W-△C
△W-△B
△W-△A
△W
RTTI
A::Wf
B::Wg
C::Wh

Virtual Base classes have *virtual thunks* which look up the offset to adjust the `this` pointer to the correct value in the virtual table!

# Compiler and Runtime Collaboration

Compiler generates:

1. . . . one code block for each method
2. . . . one virtual table for each class-composition, with
   - references to the most recent implementations of methods of a *unique common signature* (⤳ single dispatching)
   - sub-tables for the composed subclasses
   - static top-of-object and virtual bases offsets per sub-table
   - (virtual) thunks as `this`-adapters per method and subclass if needed

Runtime:

1. At program startup virtual tables are globally created
2. Allocation of memory space for each object followed by constructor calls
3. Constructor stores pointers to virtual table (or fragments) in the objects
4. Method calls transparently call methods statically or from virtual tables, *unaware of real class identity*
5. Dynamic casts may use *offset-to-top* field in objects

# Polemics of Multiple Inheritance

## Full Multiple Inheritance (FMI)

- Removes constraints on parents in inheritance
- More convenient and simple in the common cases
- Occurance of diamond pattern not as frequent as discussions indicate

## Multiple Interface Inheritance (MII)

- simpler implementation
- Interfaces and aggregation already quite expressive
- Too frequent use of FMI considered as flaw in the class hierarchy design

# Lessons Learned

**Lessons Learned**

1. Different purposes of inheritance
2. Heap Layouts of hierarchically constructed objects in C++
3. Virtual Table layout
4. LLVM IR representation of object access code
5. Linearization as alternative to explicit disambiguation
6. Pitfalls of Multiple Inheritance

- the presented approach is implemented in GNU C++ and LLVM
- Microsoft's MS VC++ approaches multiple inheritance differently
  - splits the virtual table into several smaller tables
  - keeps a vbptr (virtual base pointer) in the object representation, pointing to the virtual base of a subclass.

# Further reading...

K. Barrett, B. Cassels, P. Haahr, D. Moon, K. Playford, and T. Withington.
A monotonic superclass linearization for dylan.
In *Object Oriented Programming Systems, Languages, and Applications*, 1996.

CodeSourcery, Compaq, EDG, HP, IBM, Intel, R. Hat, and SGI.
Itanium C++ ABI.
URL: http://www.codesourcery.com/public/cxx-abi.

R. Ducournau and M. Habib.
On some algorithms for multiple inheritance in object-oriented programming.
In *Proceedings of the European Conference on Object-Oriented Programming (ECOOP)*, 1987.

R. Kleckner.
Bringing clang and llvm to visual c++ users.
URL: http://llvm.org/devmtg/2013-11/#talk11.

B. Liskov.
Keynote address – data abstraction and hierarchy.
In *Addendum to the proceedings on Object-oriented programming systems, languages and applications*, OOPSLA '87, pages 17–34, 1987.

L. L. R. Manual.
Llvm project.
URL: http://llvm.org/docs/LangRef.html.

R. C. Martin.
The liskov substitution principle.
In *C++ Report*, 1996.

P. Sabanal and M. Yason.
Reversing c++.
In *Black Hat DC*, 2007.
URL: https://www.blackhat.com/presentations/bh-dc-07/Sabanal_Yason/Paper/bh-dc-07-Sabanal_Yason-WP.pdf.

B. Stroustrup.
Multiple inheritance for C++.
In *Computing Systems*, 1999.

# Mini Seminars

1. SC=CC in Multicore Architectures with Cache (Meixner/Sorin 2006/2009)
2. Litmus Testing Memory Models: Herdtools 7
3. The Linux Kernel Memory Model
4. A Formal Analysis of the NVIDIA PTX Memory Consistency Model (2019)
5. GPU Concurrency: Weak Behaviours and Programming Assumptions (2015)
6. Transactional Memory Systems other than TSX: IBM Power 8 / BlueGene / zEnterprise
7. Lambda Calculus: Y Combinator and Recursion / SKI Combinator
8. Templates vs. Inheritance

# **Programming Languages**

Mixins and Traits

Dr. Michael Petter
Winter 2019/20

**What modularization techiques are there besides multiple implementation inheritance?**

# Outline

- Codesharing in Object Oriented Systems is often inheritance-centric
- Inheritance itself comes in different flavours:
  - single inheritance
  - multiple inheritance
- All flavours of inheritance tackle problems of *decomposition* and *composition*

# The Adventure Game

# The Adventure Game



**Door**

**<interface>Doorlike**
canPass(Person p)
canOpen(Person p)

**Short**
canPass(Person p)

**Locked**
canOpen(Person p)

**ShortLockedDoor**
canOpen(Person p)
canPass(Person p)

⚠ **Aggregation & S.-Inheritance**

- Door must explicitly provide chaining
- Doorlike must anticipate wrappers

⇒ **Multiple Inheritance** ✓

**FileStream**

read()
write()

**SocketStream**

read()
write()

?

**SynchRW**

acquireLock()
releaseLock()

**⚠ Unclear relations**

↝ Cannot inherit from both in turn with Multiple Inheritance
(*Many-to-One* instead of *One-to-Many* Relation)

# The Wrapper – Aggregation Solution



**Stream**
read()
write()

**FileStream**
read()
write()

**SocketStream**
read()
write()

**SynchRW**
read()
write()
acquireLock()
releaseLock()

⚠ **Aggregation**

- Undoes specialization
- Needs common ancestor

# The Wrapper – Multiple Inheritance Solution



---

⚠ **Duplication**

With multiple inheritance, `read`/`write` Code is essentially *identical but duplicated for each particular wrapper*

# Fragility



**⚠ Inappropriate Hierarchies**

Implemented methods (`acquireLock`/`releaseLock`) *to high*

All the problems of

- Relation
- Duplication
- Hierarchy

are centered around the question

"How do I distribute functionality over a hierarchy"

*⤳ functional (de-)composition*

# **Classes and Methods**

The building blocks for classes are

- a countable set of method *names* $\mathcal{N}$
- a countable set of method *bodies* $\mathbb{B}$

Classes map names to elements from the *flat lattice* $\mathcal{B}$ (called bindings), consisting of:

- method bodies $\in \mathbb{B}$ or classes $\in \mathcal{C}$
- $\bot$ *abstract*
- $\top$ *in conflict*

and the partial order $\bot \sqsubseteq b \sqsubseteq \top$ for each $b \in \mathcal{B}$



**Definition (Abstract Class $\in \mathcal{C}$)**

A general function $c : \mathcal{N} \mapsto \mathcal{B}$ is called a class.

**Definition (Interface and Class)**

A class $c$ is called       (with pre beeing the preimage)

> *interface* iff $\forall_{n \in \mathsf{pre}(c)} . c(n) = \bot$.
>
> *abstract class* iff $\exists_{n \in \mathsf{pre}(c)} . c(n) = \bot$.
>
> *concrete class* iff $\forall_{n \in \mathsf{pre}(c)} . \bot \sqsubset c(n) \sqsubset \top$.

# Computing with Classes and Methods

### Definition (Family of classes $\mathcal{C}$)

We call the set of all maps from names to bindings the family of classes $\mathcal{C} := \mathcal{N} \mapsto \mathcal{B}$.

Several possibilites for composing maps $\mathcal{C} \ \Box \ \mathcal{C}$:

- the symmetric join $\sqcup$, defined componentwise:

$$(c_1 \sqcup c_2)(n) = b_1 \sqcup b_2 = \begin{cases} b_2 & \text{if } b_1 = \bot \text{ or } n \notin \mathsf{pre}(c_1) \\ b_1 & \text{if } b_2 = \bot \text{ or } n \notin \mathsf{pre}(c_2) \\ b_2 & \text{if } b_1 = b_2 \\ \top & \text{otherwise} \end{cases} \quad \text{where } b_i = c_i(n)$$

- in contrast, the asymmetric join $\mathbin{\uparrow\!\!\sqcup}$, defined componentwise:

$$(c_1 \mathbin{\uparrow\!\!\sqcup} c_2)(n) = \begin{cases} c_1(n) & \text{if } n \in \mathsf{pre}(c_1) \\ c_2(n) & \text{otherwise} \end{cases}$$

# Example: Smalltalk-Inheritance

*Smalltalk* inheritance
- children's methods dominate parents' methods
- is the archetype for inheritance in mainstream languages like Java or C#
- inheriting smalltalk-style establishes a reference to the parent

**Definition (Smalltalk inheritance ($\triangleright$))**

Smalltalk inheritance is the binary operator $\triangleright : \mathcal{C} \times \mathcal{C} \mapsto \mathcal{C}$, definied by
$c_1 \triangleright c_2 = \{\texttt{super} \mapsto c_2\} \uplus (c_1 \uplus c_2)$

**Example:** Doors

$$Door = \{canPass \mapsto \bot, canOpen \mapsto \bot\}$$

$$LockedDoor = \{canOpen \mapsto 0x4204711\} \triangleright Door$$

$$= \{\texttt{super} \mapsto Door\} \uplus (\{canOpen \mapsto 0x4204711\} \uplus Door)$$

$$= \{\texttt{super} \mapsto Door, canOpen \mapsto 0x4204711, canPass \mapsto \bot\}$$

# Excursion: Beta-Inheritance

In *Beta*-style inheritance

- the design goal is to provide security wrt. replacement of a method by a different method.
- methods in parents dominate methods in subclass
- the keyword `inner` explicitly delegates control to the subclass

**Definition (Beta inheritance ($\triangleleft$))**

Beta inheritance is the binary operator $\triangleleft : \mathcal{C} \times \mathcal{C} \mapsto \mathcal{C}$, definied by
$c_1 \triangleleft c_2 = \{\texttt{inner} \mapsto c_1\} \, \mathbb{U} \, (c_2 \, \mathbb{U} \, c_1)$

Example (equivalent syntax):

```
class Person {
  String name ="Axel Simon";
  public String toString(){ return name+inner.toString();};
};
class Graduate extends Person {
  public extension String toString(){ return ", Ph.D."; };
};
```

**So what do we really want?**

# Adventure Game with Mixins

# Adventure Game with Mixins

```
class Door {
 boolean canOpen(Person p) { return true; };
 boolean canPass(Person p) { return p.size() < 210; };
}
mixin Locked {
 boolean canOpen(Person p){
  if (!p.hasItem(key)) return false; else return super.canOpen(p);
 }
}
mixin Short {
 boolean canPass(Person p){
  if (p.height()>1) return false; else return super.canPass(p);
 }
}
class ShortDoor = Short(Door);
class LockedDoor = Locked(Door);
mixin ShortLocked = Short o Locked;
class ShortLockedDoor  = Short(Locked(Door));
class ShortLockedDoor2 = ShortLocked(Door);
```

**Back to the blackboard!**

# Abstract model for Mixins

A Mixin is a *unary second order type expression*. In principle it is a curried version of the Smalltalk-style inheritance operator. In certain languages, programmers can create such mixin operators:

---

**Definition (Mixin)**

The mixin constructor $mixin : \mathcal{C} \mapsto (\mathcal{C} \mapsto \mathcal{C})$ is a unary class function, creating a unary class operator, defined by:

$$mixin(c) = \lambda x \, . \, c \triangleright x$$

---

⚠ Note: Mixins can also be composed ∘:

---

**Example:** `Doors`

$$Locked = \{canOpen \mapsto 0x1234\}$$

$$Short = \{canPass \mapsto 0x4711\}$$

$$Composed = mixin(Short) \circ (mixin(Locked)) = \lambda x \, . \, Short \, \triangleright \, (Locked \, \triangleright \, x)$$

$$= \lambda x \, . \, \{\mathtt{super} \mapsto (Locked \, \triangleright \, x)\} \, \overline{\sqcup} \, (\{canOpen \mapsto 0x1234, canPass \mapsto 0x4711\} \triangleright \, x)$$

# Wrapper with Mixins



**Mixins for wrappers**

- avoids duplication of `read`/`write` code
- keeps specialization
- even compatible to single inheritance systems

# Mixins on Implementation Level

```
class Door {
 boolean canOpen(Person p)...
 boolean canPass(Person p)...
}
mixin Locked {
 boolean canOpen(Person p)...
}
mixin Short {
 boolean canPass(Person p)...
}
class ShortDoor
   = Short(Door);
class ShortLockedDoor
   = Short(Locked(Door));
...
ShortDoor d
   = new ShortLockedDoor();
```



⚠ *non-static* super-References

⤳ dynamic dispatching without precomputed virtual table

**Surely multiple inheritance is powerful enough to simulate mixins?**

# Simulating Mixins in C++

```cpp
template <class Super>
class SyncRW : public Super {
  public: virtual int read(){
    acquireLock();
    int result = Super::read();
    releaseLock();
    return result;
  };
  virtual void write(int n){
    acquireLock();
    Super::write(n);
    releaseLock();
  };
  // ... acquireLock & releaseLock
};
```

# Simulating Mixins in C++

```cpp
template <class Super>
class LogOpenClose : public Super {
   public: virtual void open(){
    Super::open();
    log("opened");
   };
   virtual void close(){
    Super::close();
    log("closed");
   };
   protected: virtual void log(char*s) { ... };
};
class MyDocument : public SyncRW<LogOpenClose<Document>> {};
```

# True Mixins vs. C++ Mixins

## True Mixins

- `super` natively supported
- Composable mixins
- Hassle-free simple alternative to multiple inheritance

## C++ Mixins

- Mixins reduced to templated superclasses
- Can be seen as coding pattern
- C++ Type system not modular
- ⤳ Mixins have to stay source code

## Common properties of Mixins

- Linearization is necessary
- ⤳ Exact sequence of Mixins is relevant

**Ok, ok, show me a language with native mixins!**

## Ruby

```ruby
class Person
  attr_accessor :size
  def initialize
    @size = 160
  end
  def hasKey
    true
  end
end
class Door
  def canOpen (p)
    true
  end
  def canPass(person)
    person.size < 210
  end
end
```

```ruby
module Short
  def canPass(p)
    p.size < 160 and super(p)
  end
end
module Locked
  def canOpen(p)
    p.hasKey() and super(p)
  end
end

class ShortLockedDoor < Door
  include Short
  include Locked
end

p = Person.new
d = ShortLockedDoor.new
puts d.canPass(p)
```

## Ruby

```ruby
class Door
  def canOpen (p)
    true
  end
  def canPass(person)
    person.size < 210
  end
end
module Short
  def canPass(p)
    p.size < 160 and super(p)
  end
end
module Locked
  def canOpen(p)
    p.hasKey() and super(p)
  end
end
```

```ruby
module ShortLocked
  include Short
  include Locked
end
class Person
  attr_accessor :size
  def initialize
    @size = 160
  end
  def hasKey
    true
  end
end

p = Person.new
d = Door.new
d.extend ShortLocked

puts d.canPass(p)
```

**Is Inheritance the Ultimate Principle in Reusability?**

# Lack of Control



**⚠ Control**

- Common base classes are shared or duplicated at class level

# Lack of Control



**⚠ Control**

- Common base classes are shared or duplicated at class level
- `super` as ancestor reference vs. qualified specification
- ⤳ No *fine-grained specification* of duplication or sharing

# Inappropriate Hierachies



**LinkedList**

add(int, Object)
remove(int)
clear()

**Stack**

stackpointer: int

push(Object)
pushMany(Object...)
pop()

> ⚠ **Inappropriate Hierarchies**
>
> - High up specified methods *turn obsolete*, but there is no statically safe way to remove them

# Inappropriate Hierachies



```
LinkedList
─────────────────
add(int, Object)
remove(int)
clear()
```

```
Stack
─────────────────
stackpointer: int
─────────────────
push(Object)
pushMany(Object...)
pop()
```

⚠ **Inappropriate Hierarchies**

- High up specified methods *turn obsolete*, but there is no statically safe way to remove them
- ⚠ Liskov Substitution Principle!

Is Implementation Inheritance even an *Anti-Pattern*?

Excerpt from the Java 8 API documentation for class `Properties`:

> *"Because `Properties` inherits from `Hashtable`, the `put` and `putAll` methods can be applied to a `Properties` object. Their use is strongly discouraged as they allow the caller to insert entries whose keys or values are not `Strings`. The `setProperty` method should be used instead. If the `store` or `save` method is called on a "compromised" `Properties` object that contains a non-`String` key or value, the call will fail. . ."*

Excerpt from the Java 8 API documentation for class `Properties`:

> *"Because `Properties` inherits from `Hashtable`, the `put` and `putAll` methods can be applied to a `Properties` object. Their use is strongly discouraged as they allow the caller to insert entries whose keys or values are not `Strings`. The `setProperty` method should be used instead. If the `store` or `save` method is called on a "compromised" `Properties` object that contains a non-`String` key or value, the call will fail…"*

### ⚠ Misuse of Implementation Inheritance

Implementation Inheritance itself as a pattern for code reusage is often misused!
⤳ All that is not explicitly prohibited will eventually be done!

# The Idea Behind Traits

TUM

- A lot of the problems originate from the coupling of implementation and modelling
- Interfaces seem to be hierarchical
- Functionality seems to be modular

### ⚠ Central idea

Separate object *creation* from *modelling* hierarchies and *composing* functionality.

⤳ Use interfaces to design hierarchical signature propagation

⤳ Use *traits* as modules for assembling functionality

⤳ Use classes as frames for entities, which can create objects

## Traits – Composition

### Definition (Trait $\in \mathcal{T}$)

A class $t$ is without attributes is called *trait*.

The *trait sum* $+ : \mathcal{T} \times \mathcal{T} \mapsto \mathcal{T}$ is the componentwise least upper bound:
$$(c_1 + c_2)(n) = b_1 \sqcup b_2 = \begin{cases} b_2 & \text{if } b_1 = \bot \vee n \notin \mathsf{pre}(c_1) \\ b_1 & \text{if } b_2 = \bot \vee n \notin \mathsf{pre}(c_2) \\ b_2 & \text{if } b_1 = b_2 \\ \top & \text{otherwise} \end{cases} \quad \text{with } b_i = c_i(n)$$

*Trait-Expressions* also comprise:

- *exclusion* $- : \mathcal{T} \times \mathcal{N} \mapsto \mathcal{T}$: $\quad (t - a)(n) = \begin{cases} \mathsf{undef} & \text{if } a = n \\ t(n) & \text{otherwise} \end{cases}$

- *aliasing* $[\rightarrow] : \mathcal{T} \times \mathcal{N} \times \mathcal{N} \mapsto \mathcal{T}$: $\quad t[a \rightarrow b](n) = \begin{cases} t(n) & \text{if } n \neq a \\ t(b) & \text{if } n = a \end{cases}$

Traits $t$ can be connected to classes $c$ by the asymmetric join:
$$(c \mathbin{\text{\reflectbox{$\sqcap$}}} t)(n) = \begin{cases} c(n) & \text{if } n \in \mathsf{pre}(c) \\ t(n) & \text{otherwise} \end{cases}$$

Usually, this connection is reserved for the last composition level.

# Traits – Concepts

## Trait composition principles

**Flat ordering** All traits have the same precedence under $+$
$\rightsquigarrow$ explicit disambiguation with aliasing and exclusion

**Precedence** Under asymmetric join $\overline{\cup}$, class methods take precedence over trait methods

**Flattening** After asymmetric join $\overline{\cup}$: Non-overridden trait methods have the same semantics as class methods

## ⚠ Conflicts . . .

arise if composed traits map methods with identical names to different bodies

## Conflict treatment

✓ Methods can be aliased ($\rightarrow$)

✓ Methods can be excluded ($-$)

✓ Class methods override trait methods and sort out conflicts ($\overline{\cup}$)

**Can we augment classical languages by traits?**

# Extension Methods (C#)

**Central Idea:**

Uncouple method definitions from class bodies.

Purpose:

- retrospectively add methods to complex types
  $\rightsquigarrow$ *external definition*
- especially provide definitions of *interface methods*
  $\rightsquigarrow$ poor man's multiple inheritance!

**Syntax:**

1. Declare a static class with definitions of static methods
2. Explicitly declare first parameter as receiver with modifier `this`
3. Import the carrier class into scope (if needed)
4. Call extension method in *infix form* with emphasis on the receiver

```
public class Person{
 public int size = 160;
 public bool hasKey() { return true;}
}
public interface Short {}
public interface Locked {}
public static class DoorExtensions {
 public static bool canOpen(this Locked leftHand, Person p){
  return p.hasKey();
 }
 public static bool canPass(this Short leftHand, Person p){
  return p.size<160;
 }
}
public class ShortLockedDoor : Locked,Short {
 public static void Main() {
  ShortLockedDoor d = new ShortLockedDoor();
  Console.WriteLine(d.canOpen(new Person()));
 }
}
```

# Extension Methods as Traits



**Extension Methods**

- transparently extend arbitrary types externally
- provide quick relief for plagued programmers

**. . . but not traits**

- Interface declarations empty, thus kind of purposeless
- Flattening not implemented
- Static scope only

Static scope of extension methods causes unexpected errors:

```
public interface Locked {
  public bool canOpen(Person p);
}
public static class DoorExtensions {
 public static bool canOpen(this Locked leftHand, Person p){
  return p.hasKey();
 }
}
```

# Extension Methods as Traits

## Extension Methods

- transparently extend arbitrary types externally
- provide quick relief for plagued programmers

## . . . but not traits

- Interface declarations empty, thus kind of purposeless
- Flattening not implemented
- Static scope only

Static scope of extension methods causes unexpected errors:

```
public interface Locked {
  public bool canOpen(Person p);
}
public static class DoorExtensions {
 public static bool canOpen(this Locked leftHand, Person p){
  return p.hasKey();
 }
}
```

# *Virtual* **Extension Methods (Java 8)**

Java 8 advances one step further:

```java
interface Door {
  boolean canOpen(Person p);
  boolean canPass(Person p);
}
interface Locked {
  default boolean canOpen(Person p) { return p.hasKey(); }
}
interface Short {
  default boolean canPass(Person p) { return p.size<160; }
}
public class ShortLockedDoor implements Short, Locked, Door {
}
```

**Implementation**

. . . consists in adding an interface phase to `invokevirtual`'s name resolution

**⚠ Precedence**

Still, default methods do not override methods from *abstract classes* when composed

# Traits as General Composition Mechanism

> ⚠ **Central Idea**
>
> Separate class generation from hierarchy specification and functional modelling
>
> ❶ model hierarchical relations with interfaces
> ❷ compose functionality with traits
> ❸ adapt functionality to interfaces and add state via glue code in classes

**Simplified multiple Inheritance without adverse effects**

**So let's do the language with real traits?!**

# Squeak

**Smalltalk**

Squeak is a smalltalk implementation, extended with a system for traits.

Syntax:

- `name: param1 and: param2`

  declares method `name` with `param1` and `param2`
- `| ident1 ident2 |`

  declares Variables `ident1` and `ident2`
- `ident := expr`

  assignment
- `object name:content`

  sends message `name` with `content` to `object` ($\equiv$ call: `object.name(content)`)
- `.`

  line terminator
- `^ expr`

  return statement

# Traits in Squeak

```
Trait named: #TRStream uses: TPositionableStream
  on: aCollection
    self collection: aCollection.
    self setToStart.
  next
    ^ self atEnd
      ifTrue:  [nil]
      ifFalse: [self collection at: self nextPosition].
Trait named: #TSynch uses: {}
  acquireLock
    self semaphore wait.
  releaseLock
    self semaphore signal.

Trait named: #TSyncRStream uses: TSynch+(TRStream@(#readNext -> #next))
  next
    | read |
    self acquireLock.
    read := self readNext.
    self releaseLock.
    ^ read.
```

# Disambiguation

## Traits vs. Mixins vs. Class-Inheritance

All different kinds of type expressions:

- Definition of curried *second order type operators* + Linearization

*Explicitly:* Traits differ from Mixins

- Traits are applied to a class *in parallel*, Mixins *sequentially*
- Trait *composition is unordered*, avoiding linearization effects
- Traits do *not contain attributes*, avoiding state conflicts
- With traits, *glue code* is concentrated in single classes

# Disambiguation

TUM

## Traits vs. Mixins vs. Class-Inheritance

All different kinds of type expressions:

- Definition of curried *second order type operators* + Linearization
- Finegrained flat-ordered *composition of modules*

*Explicitly:* Traits differ from Mixins

- Traits are applied to a class *in parallel*, Mixins *sequentially*
- Trait *composition is unordered*, avoiding linearization effects
- Traits do *not contain attributes*, avoiding state conflicts
- With traits, *glue code* is concentrated in single classes

# Disambiguation

## Traits vs. Mixins vs. Class-Inheritance

All different kinds of type expressions:

- Definition of curried *second order type operators* + Linearization
- Finegrained flat-ordered *composition of modules*
- Definition of (local) partial order on precedence of types wrt. MRO

*Explicitly:* Traits differ from Mixins

- Traits are applied to a class *in parallel*, Mixins *sequentially*
- Trait *composition is unordered*, avoiding linearization effects
- Traits do *not contain attributes*, avoiding state conflicts
- With traits, *glue code* is concentrated in single classes

# Disambiguation

TUM

## Traits vs. Mixins vs. Class-Inheritance

All different kinds of type expressions:

- Definition of curried *second order type operators* + Linearization
- Finegrained flat-ordered *composition of modules*
- Definition of (local) partial order on precedence of types wrt. MRO
- Combination of principles

*Explicitly:* Traits differ from Mixins

- Traits are applied to a class *in parallel*, Mixins *sequentially*
- Trait *composition is unordered*, avoiding linearization effects
- Traits do *not contain attributes*, avoiding state conflicts
- With traits, *glue code* is concentrated in single classes

# Lessons learned

## Mixins

- Mixins as *low-effort* alternative to multiple inheritance
- Mixins lift type expressions to *second order type expressions*

## Traits

- Implementation Inheritance based approaches leave room for improvement in modularity in real world situations
- Traits offer *fine-grained control* of composition of functionality
- Native trait languages offer *separation of composition* of functionality from *specification* of interfaces

# Further reading...

Gilad Bracha and William Cook.
Mixin-based inheritance.
*European conference on object-oriented programming on Object-oriented programming systems, languages, and applications (OOPSLA/ECOOP)*, 1990.

James Britt.
Ruby 2.1.5 core reference, December 2014.
URL https://www.ruby-lang.org/en/documentation/.

Stéphane Ducasse, Oscar Nierstrasz, Nathanael Schärli, Roel Wuyts, and Andrew P. Black.
Traits: A mechanism for fine-grained reuse.
*ACM Transactions on Programming Languages and Systems (TOPLAS)*, 2006.

Matthew Flatt, Shriram Krishnamurthi, and Matthias Felleisen.
Classes and mixins.
*Principles of Programming Languages (POPL)*, 1998.

Brian Goetz.
Interface evolution via virtual extension methods.
*JSR 335: Lambda Expressions for the Java Programming Language*, 2011.

Anders Hejlsberg, Scott Wiltamuth, and Peter Golde.
*C# Language Specification*.
Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2003.
ISBN 0321154916.

Nathanael Schärli, Stéphane Ducasse, Oscar Nierstrasz, and Andrew P. Black.
Traits: Composable units of behaviour.
*European Conference on Object-Oriented Programming (ECOOP)*, 2003.

# **Programming Languages**

Prototypes

Dr. Michael Petter
Winter 2019/20

**Prototype based programming**

1. Basic language features
2. Structured data
3. Code reusage
4. Imitating Object Orientation

"Why bother with modelling types for my quick hack?"

## Bothersome features

- Specifying types for singletons
- Getting generic types right inspite of co- and contra-variance
- Subjugate language-imposed inheritance to (mostly) avoid redundancy

## Prototype based programming

- Start by creating examples
- Only very basic concepts
- Introduce complexity only by need
- Shape language features yourself!

"Let's go back to basic concepts – *Lua*"

# Basic Language Features

- Chunks being sequences of statements.
- Global variables implicitely defined

```
s = 0;
i = 1              -- Single line comment
p = i+s p=42       --[[ Multiline
comment --]]
s = l
```

# Basic Types and Values

- Dynamical types – no type definitions
- Each value carries its type
- `type()` returns a string representation of a value's type

```lua
a = true
type(a)              -- boolean
type("42"+0)         -- number
type("Petter "..1)   -- string
type(type)           -- function
type(nil)            -- nil
type([[<html><body>pretty long string</body>
</html>
]])                  -- string
a = 42
type(a)              -- number
```

# Functions for Code

✓ First class citizens

```
function prettyprint(title, name, age)
  return title.." "..name..", born in "..(2018-age)
end

a = prettyprint
a("Dr.","Petter",42)

prettyprint = function (title, name, age)
  return name..", "..title
end
```

# Introducing Structure

- only one complex data type
- indexing via arbitrary values *except nil* ($\rightsquigarrow$ Runtime Error)
- arbitrary large and dynamically growing/shrinking

```lua
a = {}                 -- create empty table
k = 42
a[k] = 3.14159         -- entry 3.14159 at key 42
a["k"] = k             -- entry 42 at key "k"
a[k] = nil             -- deleted entry at key 42
print(a.k)             -- syntactic sugar for a["k"]
```

## Table Lifecycle

- created from scratch
- modification is persistent
- assignment with reference-semantics
- garbage collection

```
a = {}                  -- create empty table
a.k = 42
b = a                   -- b refers to same as a
b["k"] = "k"            -- entry "k" at key "k"
print(a.k)              -- yields "k"
a = nil
print(b.k)              -- still "k"
b = nil
print(b.k)              -- nil now
```

"So far nothing special – let's compose types"

# Table Behaviour

## Metatables

- are *ordinary tables*, used as collections of special functions
- Naming conventions for special functions
- Connect to a table via `setmetatable`, retrieve via `getmetatable`
- Changes behaviour of tables

```lua
meta = {}                          -- create as plain empty table
function meta.__tostring(person)
  return person.prefix .. " " .. person.name
end
a = { prefix="Dr.",name="Petter"} -- create Michael
setmetatable(a,meta)               -- install metatable for a
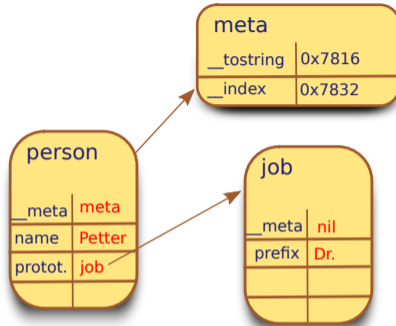print(a)                           -- print "Dr. Petter"
```

- Overload operators like `__add,__mul,__sub,__div,__pow,__concat,__unm`
- Overload comparators like `__eq,__lt,__le`

## Delegation

- ⚠ reserved key `__index` determines *handling* of failed name lookups
- convention for signature: receiver table and key as parameters
- if dispatching to another table ↝ *Delegation*

```
meta = {}
function meta.__tostring(person)
  return person.prefix .. " " .. person.name
end
function meta.__index(tbl, key)
  return tbl.prototype[key]
end
job = { prefix="Dr." }
person = { name="Petter",prototype=job } -- create Michael
setmetatable(person,meta)               -- install metatable
print(person)                           -- print "Dr. Petter"
```

```
function meta.__tostring(person)  -- 0x7816
  return person.prefix .. " " .. person.name
end
function meta.__index(tbl, key) -- 0x7832
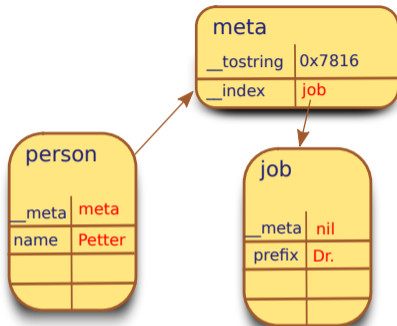  return tbl.prototype[key]
end
```

# Delegation 2

⤳ Conveniently, `__index` does not need to be a function

```lua
meta = {}
function meta.__tostring(person)
  return person.prefix .. " " .. person.name
end
job = { prefix="Dr." }
meta.__index = job                    -- delegate to job
person = { name="Petter" }            -- create Michael
setmetatable(person,meta)             -- install metatable
print(person)                         -- print "Dr. Petter"
```

```lua
function meta.__tostring(person)  -- 0x7816
  return person.prefix .. " " .. person.name
end
```

## Delegation 3

- `__newindex` handles unresolved updates
- frequently used to implement protection of objects

```
meta = {}
function meta.__newindex(abl,key,val)
  if (key == "title" and tbl.name=="Guttenberg") then
    error("No title for You, sir!")
  else
    tbl.data[key]=val
  end
end
function meta.__tostring(tbl)
  return (tbl.title or "") .. table.name
end
person={ data={} }            -- create person's data
meta.__index = person.data
setmetatable(person,meta)
person.name = "Guttenberg"    -- name KT
person.title = "Dr."          -- try to give him Dr.
```

# Object Oriented Programming

⚠ so far no concept for multiple *objects*

```
Account = { balance=0 }
function Account.withdraw (val)
 Account.balance=Account.balance-val
end
function Account.__tostring()
   return "Balance is "..Account.balance
end
setmetatable(Account,Account)
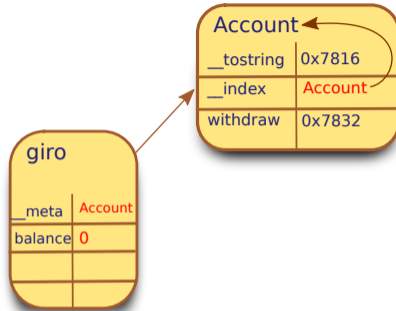Account.withdraw(10)
print(Account)
```

# Introducing Identity

- Concept of an object's *own identity* via parameter
- Programming aware of multiple instances
- Share code between instances

```lua
function Account.withdraw (acc, val)
 acc.balance=acc.balance-val
end
function Account.tostring(acc)
   return "Balance is "..acc.balance
end
Account.__index=Account        -- share Account's functions
mikes = { balance = 0 }
daves = { balance = 0 }
setmetatable(mikes,Account)     -- delegate from mikes to Account
setmetatable(daves,Account)     -- del. from daves to Account
Account.withdraw(mikes,10)
mikes.withdraw(mikes,10)        -- withdraw independently
mikes:withdraw(10)
print(daves:tostring() .. " " .. mikes:tostring())
```

```
function Account.withdraw (acc, val)
 acc.balance=acc.balance-val
end
function Account.tostring(acc)
   return "Balance is "..acc.balance
end
```

# Introducing "Classes"

- Particular tables *used* like classes
- *self* table for accessing object-relative attributes
- connection via creator function *new* (like a constructor)

```
function Account:withdraw (val)
 self.balance=self.balance-val
end
function Account:tostring()
    return "Balance is "..self.balance
end
function Account:new(template)
 template = template or {balance=0}   -- initialize
 setmetatable(template,{__index=self})-- delegate to Account
 getmetatable(template).__tostring = Account.tostring
 return template
end
giro = Account:new({balance=10})      -- create instance
giro:withdraw(10)
print(giro)
```

# Inheriting Functionality

- Differential description possible in child class style
- Easily creating particular singletons

```
LimitedAccount = { }
setmetatable(LimitedAccount,{__index=Account})
function LimitedAccount:new()
  instance = { balance=0,limit=100 }
  setmetatable(instance,{__index=self})
end
function LimitedAccount:withdraw(val)
  if (self.balance+self.limit < val) then
     error("Limit exceeded")
  end
  Account.withdraw(self,val)
end
specialgiro = LimitedAccount:new()
specialgiro:withdraw(90)
print(specialgiro)
```

## Multiple Inheritance

⤳ Delegation leads to chain-like inheritance

```lua
function createClass (parent1,parent2)
  local c = {}                    -- new class, child of p1&p2
  setmetatable(c, {__index =
    function (t, k)               -- search for each name
      local v = parent1[k]        -- in both parents
      if v then return v end
      return parent2[k]
    end}
  )
  c.__index = c                   -- c is prototype of instances
  function c:new (o)              -- constructor for this class
    o = o or {}
    setmetatable(o, c)            -- c is also metatable
    return o
  end
  return c                        -- finally return c
end
```

```
Doctor    = { postfix="Dr. "}
Researcher = { prefix=" ,Ph.D."}

ResearchingDoctor = createClass(Doctor,Researcher)
axel = ResearchingDoctor:new( { name="Michael Petter" } )
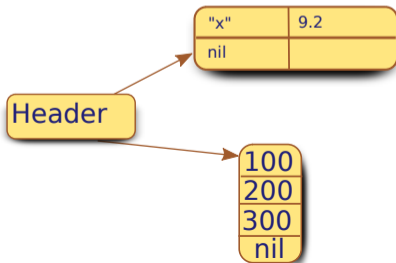print(axel.prefix..axel.name..axel.postfix)
```

↝ The special case of dual-inheritance can be extended to comprise multiple inheritance

# Implementation of Lua

```
typedef struct {
  int type_id;
  Value v;
} TObject;
```

```
typedef union {
  void *p;
  int b;
  lua_number n;
  GCObject *gc;
} Value;
```

- Datatypes are simple values (Type+union of different flavours)
- Tables at low-level fork into Hashmaps with pairs and an integer-indexed array part

# Further Topics in Lua

- Coroutines
- Closures
- Bytecode & Lua-VM

## Lessons Learned

1. Abandoning fixed inheritance yields ease/speed in development
2. Also leads to horrible runtime errors
3. Object-orientation and multiple-inheritance as special cases of delegation
4. Minimal featureset eases implementation of compiler/interpreter
5. Room for static analyses to find bugs ahead of time

**Further Reading...**

📕 Roberto Ierusalimschy.
*Programming in Lua, Third Edition*.
Lua.Org, 2013.
ISBN 859037985X.

📕 Roberto Ierusalimschy, Luiz Henrique de Figueiredo, and Waldemar Celes Filho.
Lua-an extensible extension language.
*Softw., Pract. Exper.*, 1996.

📕 Roberto Ierusalimschy, Luiz Henrique de Figueiredo, and Waldemar Celes.
The implementation of lua 5.0.
*Journal of Universal Computer Science*, 2005.

# **Programming Languages**

Aspect Oriented Programming

Dr. Michael Petter
Winter 2019/20

"Is modularity the key principle to organizing software?"

## Learning outcomes

1. AOP Motivation and Weaving basics
2. Bundling aspects with static crosscutting
3. Join points, Pointcuts and Advice
4. Composing Pointcut Designators
5. Implementation of Advices and Pointcuts

# **Motivation**

- Traditional modules directly correspond to code blocks
- Aspects can be thought of seperately but are smeared over modules ⤳ *Tangling of aspects*
- Focus on *Aspects of Concern*

⤳ *A*spect *O*riented *P*rogramming

# Motivation

- Traditional modules directly correspond to code blocks
- Aspects can be thought of seperately but are smeared over modules ⤳ *Tangling of aspects*
- Focus on *Aspects of Concern*

⤳ *A*spect *O*riented *P*rogramming

## Aspect Oriented Programming

- Express a system's aspects of concerns cross-cutting modules
- Automatically combine separate Aspects with a *Weaver* into a program

Functional decomposition

Compiler

Functional decomposition

Aspect oriented decomposition

Compiler

Aspect Weaver

# System Decomposition in Aspects

Example concerns:

- Security
- Logging
- Error Handling
- Validation
- Profiling

# System Decomposition in Aspects

Example concerns:

- Security
- Logging
- Error Handling
- Validation
- Profiling

$\rightsquigarrow$     AspectJ

# Static Crosscutting

# **Adding External Defintions**

inter-type declaration

```
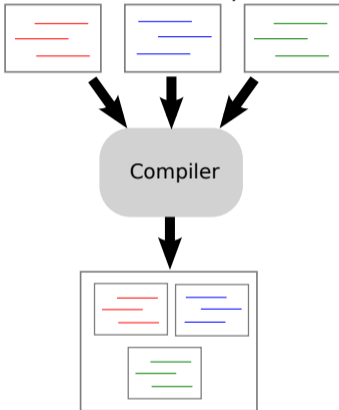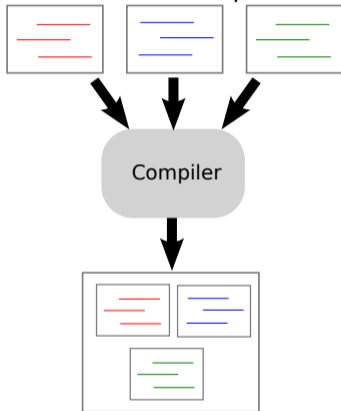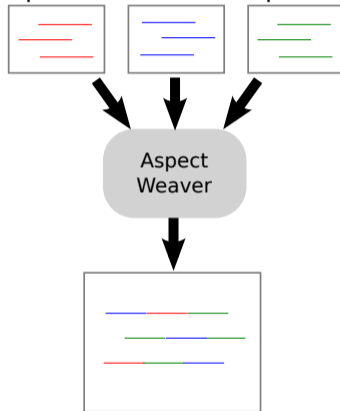class Expr {}
class Const extends Expr {
  public int val;
  public Const(int val) {
    this.val=val;
} }
class Add extends Expr {
  public Expr l,r;
  public Add(Expr l, Expr r) {
   this.l=l;this.r=r;
} }

aspect ExprEval {
  abstract int Expr.eval();
  int Const.eval(){ return val; };
  int Add.eval()  { return l.eval()
                        + r.eval(); }
}
```

⇒

equivalent code

```
// aspectj-patched code
abstract class Expr {
  abstract int eval();
}
class Const extends Expr {
  public int val;
  public int eval(){ return val; };
  public Const(int val) {
    this.val=val;
} }
class Add extends Expr {
  public Expr l,r;
  public int eval() { return l.eval()
                          + r.eval(); }
  public Add(Expr l, Expr r) {
   this.l=l;this.r=r;
  }
}
```

# Dynamic Crosscutting

# Join Points

Well-defined points in the control flow of a program

| | |
|---|---|
| method/constr. call | executing the actual method-call statement |
| method/constr. execution | the individual method is executed |
| field get | a field is read |
| field set | a field is set |
| exception handler execution | an exception handler is invoked |
| class initialization | static initializers are run |
| object initialization | dynamic initializers are run |

# Pointcuts and Designators

## Definition (Pointcut)

A pointcut is a *set of join points* and optionally some of the runtime values when program execution reaches a refered join point.

Pointcut designators can be defined and named by the programmer:

⟨*userdef*⟩ ::= 'pointcut' ⟨*id*⟩ '(' ⟨*idlist*⟩$^?$ ')' ':' ⟨*expr*⟩ ';'

⟨*idlist*⟩ ::= ⟨*id*⟩ ( ',' ⟨*id*⟩ )*

⟨*expr*⟩ ::= '!' ⟨*expr*⟩
  | ⟨*expr*⟩ '&&' ⟨*expr*⟩
  | ⟨*expr*⟩ '||' ⟨*expr*⟩
  | '(' ⟨*expr*⟩ ')'
  | ⟨*primitive*⟩

Example:

```
pointcut dfs(): execution (void Tree.dfs()) ||
                execution (void Leaf.dfs()) ;
```

## Advice

... are method-like constructs, used to define additional behaviour at joinpoints:

- before(formal)
- after(formal)
- after(formal) returning (formal)
- after(formal) throwing (formal)

For example:

```
aspect Doubler {
  before(): call(int C.foo(int)) {
    System.out.println("About to call foo");
} }
```

# Binding Pointcut Parameters in Advices

Certain pointcut primitives add dependencies on the context:

- args(arglist)

This binds identifiers to parameter values for use in in advices.

```
aspect Doubler {
  before(int i): call(int C.foo(int)) && args(i) {
      i = i*2;
} }
```

arglist actually is a flexible expression:

$\langle \text{arglist} \rangle ::= ( \langle \text{arg} \rangle ( \text{`,'} \langle \text{arg} \rangle )^* )^?$

| $\langle \text{arg} \rangle ::=$ | $\langle \text{identifier} \rangle$ | binds a value to this identifier |
| | $\langle \text{typename} \rangle$ | filters only this type |
| | `*` | matches all types |
| | `..` | matches several arguments |

# Around Advice

Unusual treatment is necessary for

- `type around(formal)`

⚠ Here, we need to pinpoint, where the advice is wrapped around the join point – this is achieved via `proceed()`:

```
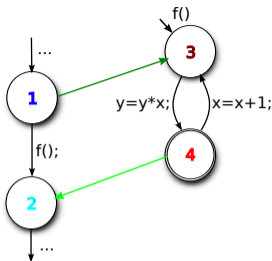aspect Doubler {
  int around(int i): call(int C.foo(Object, int)) && args(i) {
    int newi = proceed(i*2);
    return newi/2;
} }
```

# Pointcut Designator Primitives

# Method Related Designators



- `call(signature)`
- `execution(signature)`

Matches call/execution join points at which the method or constructor called matches the given *signature*. The syntax of a method/constructor *signature* is:

```
ResultTypeName RecvrTypeName.meth_id(ParamTypeName, ...)
NewObjectTypeName.new(ParamTypeName, ...)
```

# Method Related Designators

```java
class MyClass{
  public String toString() {
    return "silly me ";
  }
  public static void main(String[] args){
    MyClass c = new MyClass();
    System.out.println(c + c.toString());
} }
aspect CallAspect {
  pointcut calltostring() : call     (String MyClass.toString());
  pointcut exectostring() : execution(String MyClass.toString());
  before() :  calltostring() || exectostring() {
    System.out.println("advice!");
} }
```

# Method Related Designators

```java
class MyClass{
  public String toString() {
    return "silly me ";
  }
  public static void main(String[] args){
    MyClass c = new MyClass();
    System.out.println(c + c.toString());
} }
aspect CallAspect {
  pointcut calltostring() : call     (String MyClass.toString());
  pointcut exectostring() : execution(String MyClass.toString());
  before() :  calltostring() || exectostring() {
    System.out.println("advice!");
} }
```

```
advice!
advice!
advice!
silly me silly me
```

# Field Related Designators

- get(fieldqualifier)
- set(fieldqualifier)

Matches field get/set join points at which the field accessed matches the signature. The syntax of a field qualifier is:

```
FieldTypeName ObjectTypeName.field_id
```

⚠ : However, set has an argument which is bound via args:

```
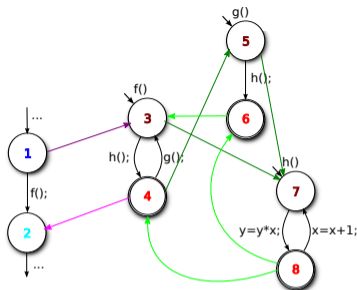aspect GuardedSetter {
  before(int newval): set(static int MyClass.x) && args(newval) {
    if (Math.abs(newval - MyClass.x) > 100)
      throw new RuntimeException();
} }
```

# Type based

- target(typeorid)
- within(typepattern)
- withincode(methodpattern)

Matches join points of any kind which

- are refering to the receiver of type typeorid
- is contained in the class body of type typepattern
- is contained within the method defined by methodpattern

# Flow and State Based



- `cflow(arbitrary_pointcut)`

Matches join points of *any kind* that occur strictly between entry and exit of each join point matched by `arbitrary_pointcut`.

- `if(boolean_expression)`

Picks join points based on a dynamic property:

```
aspect GuardedSetter {
  before(): if(thisJoinPoint.getKind().equals(METHOD_CALL)) && within(MyClass) {
    System.out.println("What an inefficient way to match calls");
} }
```

# Which advice is served first?

## Advices are defined in different aspects

- If statement `declare precedence:A, B;` exists, then advice in aspect `A` has precedence over advice in aspect `B` for the same join point.
- Otherwise, if aspect `A` is a subaspect of aspect `B`, then advice defined in `A` has precedence over advice defined in `B`.
- Otherwise, (i.e. if two pieces of advice are defined in two different aspects), it is *undefined* which one has precedence.

## Advices are defined in the same aspect

- If either are *after advice*, then the one that appears *later* in the aspect has precedence over the one that appears earlier.
- Otherwise, then the one that appears *earlier* in the aspect has precedence over the one that appears later.

**Implementation**

# Implementation

Aspect Weaving:

- Pre-processor
- During compilation
- Post-compile-processor
- During Runtime in the Virtual Machine
- A combination of the above methods

# Woven JVM Code

```
Expr one = new Const(1);
one.val = 42;
```

```
aspect MyAspect {
  pointcut settingconst(): set(int Const.val);
  before () : settingconst() {
    System.out.println("setter");
} }
```

```
...
117: aload_1
118: iconst_1
119: dup_x1
120: invokestatic  #73 // Method MyAspect.aspectOf:()LMyAspect;
123: invokevirtual #79 // Method MyAspect.ajc$before$MyAspect$2$704a2754:()V
126: putfield      #54 // Field Const.val:I
...
```

# Woven JVM Code

```
Expr one = new Const(1);
Expr e = new Add(one,one);
String s = e.toString();
System.out.println(s);
```

```
aspect MyAspect {
  pointcut callingtostring():
    call (String Object.toString()) && target(Expr);
  before () : callingtostring() {
    System.out.println("calling");
} }
```

```
...
 72: aload_2
 73: instanceof    #1  // class Expr
 76: ifeq          85
 79: invokestatic  #67 // Method MyAspect.aspectOf:()MyAspect;
 82: invokevirtual #70 // Method MyAspect.ajc$before$MyAspect$1$4c1f7c11:()V
 85: aload_2
 86: invokevirtual #33 // Method java/lang/Object.toString:()Ljava/lang/String;
 89: astore_3
...
```

# Poincut Parameters and Around/Proceed

Around clauses often refer to parameters and `proceed()` – sometimes across different contexts!

```
class C {
  int foo(int i) { return 42+i; }
}
aspect Doubler {
  int around(int i): call(int *.foo(Object, int)) && args(i) {
    int newi = proceed(i*2);
    return newi/2;
} }
```

⚠ Now, imagine code like:

```
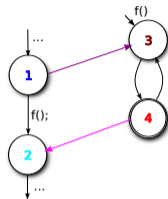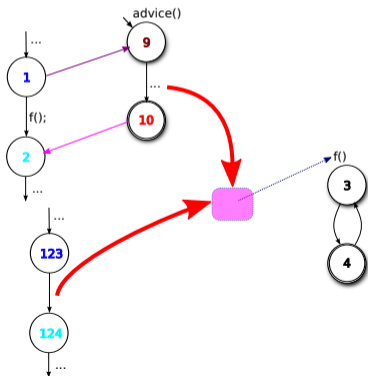public static void main(String[] args){
  new C().foo(42);
}
```

# Around/Proceed – via Procedures

✓ inlining advices in `main` – all of it in JVM, disassembled to equivalent:

```
// aspectj patched code
public static void main(String[] args){
  C c = new C();
  foo_aroundBody1Advice(c,42,Doubler.aspectOf(),42,null);
}
private static final int foo_aroundBody0(C c, int i){
  return c.foo(i);
}
private static final int foo_aroundBody1Advice
    (C c, int i, Doubler d, int j, AroundClosure a){
      int temp = 2*i;
      int ret = foo_aroundBody0(c,temp);
      return ret / 2;
}
```

⚠ However, instead of beeing used for a direct call, `proceed()` and its parameters may *escape the calling context*:

# Pointcut parameters and Scope

⚠ `proceed()` might not even be in the same scope as the original method!

⚠ even worse, the scope of the exposed parameters might have expired!

```
class C {
  int foo(int i) { return 42+i; }
  public static void main(String[] str){ new C().foo(42); }
}
aspect Doubler {
    Executor executor;
    Future<Integer> f;
    int around(int i): call(int *.foo(Object, int)) && args(i) {
      Callable<Integer> c = () -> proceed(i*2)/2;
      f = executor.submit(c);
      return i/2;
    }
    public int getCachedValue() throws Exception {
        return f.get();
    }
} }
```

# Shadow Classes and Closures



✓ creates a shadow, carrying the advice

✓ creates a closure, carrying the context/parameters

```
// aspectj patched code
public static void main(String[] str){
  int itemp = 42;
  Doubler shadow = Doubler.aspectOf();
  Object[] params = new Object[]
    { new C(),Conversions.intObject(itemp) };
  C_AjcClosure1 closure = new C_AjcClosure1(params);
  shadow.ajc$around$Doubler$1$9158ff14(itemp,closure);
}
```

# Shadow Classes and Closures

```
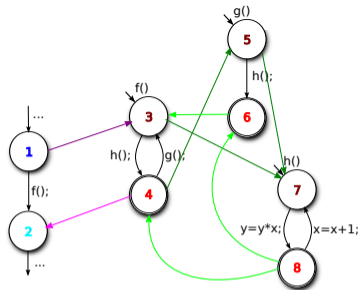// aspectj patched code
class Doubler {    // shadow class, holding the fields for the advice
  Future<Integer> f;
  ExecutorService executor;
  ...
  public int ajc$around$Doubler$1$9158ff14(int i, AroundClosure c){
    Callable<Integer> c = lambda$0(i,c);
    f = executor.submit(c);
    return i/2;
    }
  public static int ajc$around$Doubler$1$9158ff14proceed(int i, AroundClosure c)
    throws Throwable{
     Object[] params = new Object[] { Conversions.intObject(i) };
     return Conversions.intValue(c.run(params));
  }
  static Integer lambda$0(int i, AroundClosure c) throws Exception{
    return Integer.valueOf(ajc$around$Doubler$1$9158ff14proceed(i*2, c)/2);
} }
class C_AjcClosure1 extends AroundClosure{ // closure class for poincut params
  C_AjcClosure1(Object[] params){ super(params); }
  Object run(Object[] params) {
    C c = (C) params[0];
    int i = Conversions.intValue(params[1]);
    return Conversions.intObject(C.foo_aroundBody0(c, i));
} }
```

# Property Based Crosscutting



```
after(int i) : call(void h()) &&
        cflow( call(void f(int)) &&
                args(i))
        { ... } ;
```

## Idea 1: Stack based

- At each `call`-match, check runtime stack for `cflow`-match
- ⤳ Naive implementation
- ⤳ Poor runtime performance

## Idea 2: State based

- Keep seperate stack of states
- ⤳ Only modify stack at `cflow`-relevant pointcuts
- ⤳ Check stack for emptyness

Even more optimizations in practice
⤳ state-sharing, ⤳ counters,
⤳ static analysis

# Implementation – Summary

Translation scheme implications:

**before/after Advice** ... ranges from *inlined code* to distribution into *several methods and closures*

**Joinpoints** ... in the original program that have advices may get *explicitly dispatching wrappers*

**Dynamic dispatching** ... can require a *runtime test* to correctly interpret certain joinpoint designators

**Flow sensitive pointcuts** ... runtime penalty for the naive implementation, optimized version still *costly*

# Aspect Orientation

## Pro

- Un-tangling of concerns
- Late extension across boundaries of hierarchies
- Aspects provide another level of abstraction

## Contra

- Weaving generates runtime overhead
- nontransparent control flow and interactions between aspects
- Debugging and Development needs IDE Support

# Further reading...

📕 Pavel Avgustinov, Aske Simon Christensen, Laurie Hendren, Sascha Kuzins, Jennifer Lhoták, Ondřej Lhoták, Oege de Moor, Damien Sereni, Ganesh Sittampalam, and Julian Tibble.
Optimising aspectj.
*SIGPLAN Not.*, 40(6):117–128, June 2005.

📕 Gregor Kiczales.
Aspect-oriented programming.
*ACM Comput. Surv.*, 28(4es), 1996.
ISSN 0360-0300.

📕 Gregor Kiczales, Erik Hilsdale, Jim Hugunin, Mik Kersten, Jeffrey Palm, and WilliamG. Griswold.
An overview of aspectj.
*ECOOP 2001 — Object-Oriented Programming*, 2072:327–354, 2001.

📕 H. Masuhara, G. Kiczales, and C. Dutchyn.
A compilation and optimization model for aspect-oriented programs.
*Compiler Construction*, 2622:46–60, 2003.

# **Programming Languages**

Metaprogramming

Dr. Michael Petter
Winter 2019/20

"Let's write a program, which writes a program"

## Learning outcomes

1. Compilers and Compiler Tools
2. Preprocessors for syntax rewriting
3. Reflection and Metaclasses
4. Metaobject Protocol
5. Macros

## **Motivation**

- Aspect Oriented Programming establishes programmatic refinement of program code
- How about establishing support for program refinement in the language concept itself?
- Treat program *code as data*

⤳ Metaprogramming

# Motivation

- Aspect Oriented Programming establishes programmatic refinement of program code
- How about establishing support for program refinement in the language concept itself?
- Treat program *code as data*

⤳ Metaprogramming

## Metaprogramming

- Treat programs as data
- Read, analyse or transform (other) programs
- Program modifies itself during runtime

**Codegeneration Tools**

# Codegeneration Tools

## Compiler Construction

In Compiler Construction, there are a lot of codegeneration tools, that compile DSLs to target source code. Common examples are `lex` and `bison`.

Example: `lex`:
`lex` generates a table lookup based implementation of a finite automaton corresponding to the specified disjunction of regular expressions.

```
%{ #include <stdio.h>
%}
%%          /* Lexical Patterns */
[0-9]+   {    printf("integer: %s\n", yytext); }
.|\n     {    /* ignore            */              }
%%
int main(void) {
    yylex();
    return 0;
}
```

**Codegeneration via Preprocessor**

# Compiletime-Codegeneration

## String Rewriting Systems

A Text Rewriting System provides a set of grammar-like rules (→*Macros*) which are meant to be applied to the target text.

Example: *C P*re*p*rocessor (CPP)

```
#define min(X,Y) (( X < Y )? (X) : (Y))
x = min(5,x);    // (( 5 < x )? (5) : (x))
x = min(++x,y+5); // (( ++x < y+5 )? (++x) : (y+5))
```

# Compiletime-Codegeneration

## String Rewriting Systems

A Text Rewriting System provides a set of grammar-like rules (→*Macros*) which are meant to be applied to the target text.

Example: *C P*re*p*rocessor (CPP)

```
#define min(X,Y) (( X < Y )? (X) : (Y))
x = min(5,x);    // (( 5 < x )? (5) : (x))
x = min(++x,y+5); // (( ++x < y+5 )? (++x) : (y+5))
```

### ⚠ Nesting, Precedence, Binding, Side effects, Recursion, . . .

- Parts of Macro parameters can bind to context operators depending on the precedence and binding behaviour
- Side effects are recomputed for every occurance of the Macro parameter
- Any (indirect) recursive replacement stops the rewriting process
- Name spaces are not separated, identifiers duplicated

# Compiletime-Codegeneration

Example application: Language constructs [3]:

```
ATOMIC {
  i--;
  i++;
}
```

```
#define ATOMIC            \
  acquire(&globallock);\
  { /* user code */ }  \
  release(&globallock);
```

## Compiletime-Codegeneration

Example application: Language constructs [3]:

```
ATOMIC {
  i--;
  i++;
}
```

```
#define ATOMIC           \
  acquire(&globallock);\
  { /* user code */ }  \
  release(&globallock);
```

⚠ How can we bind the block, following the ATOMIC to the usercode fragment?

Particularly in a situation like this?

```
if (i>0)
  ATOMIC {
    i--;
    i++;
  }
```

# Compiletime-Codegeneration

Append code to usercode

```
if (1)
  goto body;
else
  while (1)
    if (1) {
      /* appended code */
      break;
    }
    else body:
    {/* block following the macro */}
```

Prepend code to usercode

```
if (1)
  /* prepended code */
  goto body;
else
  body:
  {/* block following the macro */}
```

# Compiletime-Codegeneration

All in one

```c
if (1) {
  /* prepended code */
  goto body;
} else
  while (1)
    if (1) {
      /* appended code */
      break;
    }
    else body:
    { /* block following the expanded macro */ }
```

# Compiletime-Codegeneration

```
#define concat_( a, b) a##b
#define label(prefix, lnum) concat_(prefix,lnum)
#define ATOMIC                  \
if (1) {                        \
  acquire(&globallock);         \
  goto label(body,__LINE__);    \
} else                          \
    while (1)                   \
      if (1) {                  \
        release(&globallock);   \
        break;                  \
      }                         \
      else                      \
        label(body,__LINE__):
```

> ⚠ **Reusability**
>
> labels have to be created dynamically in order for the macro to be reusable ($\rightarrow$ __LINE__)

**Homoiconic Metaprogramming**

# Homoiconic Programming

## Homoiconicity

In a homoiconic language, the primary representation of programs is also a data structure in a primitive type of the language itself.

**data is code**
**code is data**

- Metaclasses and Metaobject Protocol
- (Hygienic) Macros

**Reflection**

# Reflective Metaprogramming

## Type introspection

A language with *Type introspection* enables to examine the type of an object at runtime.

Example: Java `instanceof`

```java
public boolean equals(Object o){
  if (!(o instanceof Natural)) return false;
  return ((Natural)o).value == this.value;
}
```

# Reflective Metaprogramming

Metaclasses ($\rightarrow$ **code is data**)

Example: Java Reflection / Metaclass `java.lang.Class`

```java
static void fun(String param){
  Object incognito = Class.forName(param).newInstance();
  Class meta = incognito.getClass(); // obtain Metaobject
  Field[] fields = meta.getDeclaredFields();
  for(Field f : fields){
    Class t = f.getType();
    Object v = f.get(o);
    if(t == boolean.class && Boolean.FALSE.equals(v))
    // found default value
    else if(t.isPrimitive() && ((Number) v).doubleValue() == 0)
    // found default value
    else if(!t.isPrimitive() && v == null)
    // found default value
} }
```

**Metaobject Protocol**

# Metaobject Protocol

Metaobject Protocol (MOP [1])

Example: Lisp's CLOS metaobject protocol

... offers an interface to manipulate the underlying implementation of CLOS to adapt the system to the programmer's liking in aspects of

- creation of classes and objects
- creation of new properties and methods
- causing inheritance relations between classes
- creation generic method definitions
- creation of method implementations
- creation of specializers ($\rightarrow$ overwriting, multimethods)
- configuration of standard method combination ($\rightarrow$ before,after,around, call-next-method)
- simple or custom method combinators ($\rightarrow$ +,append,max,...)
- addition of documentation

**Hygienic Macros**

# Homoiconic Runtime-Metaprogramming

## Clojure! [2]

Clojure programs are represented after parsing in form of symbolic expressions (*S-Expressions*), consisting of nested trees:

### S-Expressions

S-Expressions are either

- an atom
- an expression of the form $(x.y)$ with $x, y$ being S-Expressions

Remark: Established shortcut notation for lists:

$$(x_1\ x_2\ x_3) \quad \equiv \quad (x_1\ .\ (x_2\ .\ (x_3\ .\ ())))$$

# Homoiconic Runtime-Metaprogramming

## Special Forms

Special forms differ in the way that they are interpreted by the clojure runtime from the standard evaluation rules.

Language Implementation Idea: reduce every expression to special forms:

```
(def symbol doc? init?)
(do expr*)
(if test then else?)
(let [binding*] expr*)
(eval form)  ; evaluates the datastructure form
(quote form) ; yields the unevaluated form
(var symbol)
(fn name? ([params*] expr*)+)
(loop [binding*] expr*)
(recur expr*) ; rebinds and jumps to loop or fn
;...
```

# Homoiconic Runtime-Metaprogramming

## Macros

Macros are configurable syntax/parse tree transformations.

Language Implementation Idea: define advanced language features in macros, based very few *special forms* or other macros.

Example: While loop:

```
(macroexpand '(while a b))
; => (loop* [] (clojure.core/when a b (recur)))

(macroexpand '(when a b))
;=> (if a (do b))
```

# Homoiconic Runtime-Metaprogramming

Macros can be written by the programmer in form of S-Expressions:

```
(defmacro infix
  "converting infix to prefix"
  [infixed]
  (list (second infixed) (first infixed) (last infixed)))
```

...producing

```
(infix (1 + 1))
; => 2
(macroexpand '(infix (a + b)))
; => (+ a b)
```

### ⚠ Quoting

Macros and functions are directly interpreted, if not *quoted* via

```
(quote keyword)  ; or equivalently:
'keyword
; => keyword
```

# Homoiconic Runtime-Metaprogramming

```
(defmacro fac1 [n]
   (if (= n 0)
      1
      (list '* n (list 'fac1 (- n 1)
 ))))
```

```
(defn fac2 [n]
   (if (= n 0)
      1
      (* n (fac2 (- n 1)
 ))))
```

```
(fac1 4)
; => 24
```

```
(fac2 4)
; => 24
```

...produces

```
(macroexpand '(fac1 4))
; => (* 4 (fac1 3))

(macroexpand-all '(fac1 4))
; => (* 4 (* 3 (* 2 (* 1 1))))
```

⤳ why bother?

# Homoiconic Runtime-Metaprogramming

### ⚠ Macros vs. Functions

- Macros as static AST Transformations, vs. Functions as runtime control flow manipulations
- Macros replicate parameter forms, vs. Functions evaluate parameters once
- ⤳ Macro parameters are uninterpreted, not necessarily valid expressions, vs. Functions parameters need to be valid expressions

# Homoiconic Runtime-Metaprogramming

## ⚠ Macro Hygiene

*Shadowing* of variables may be an issue in macros, and can be avoided by generated symbols!

```
(def variable 42)
(macro mac [&stufftodo] `(let [variable 4711] ~@stufftodo))
(mac (println variable))
; => can't let qualified name: variable
```

```
(macro mac [&stufftodo] `(let [variable# 4711] ~@stufftodo))
```

⤳ Symbol generation to avoid namespace collisions!

# Further reading...

📖 Richard P. Gabriel.
Gregor kiczales, jim des rivières, and daniel g. bobrow, the art of the metaobject protocol.
*Artif. Intell.*, 61(2):331–342, 1993.
URL: https://doi.org/10.1016/0004-3702(93)90073-K,
doi:10.1016/0004-3702(93)90073-K.

📖 Daniel Higginbotham.
*Clojure for the Brave and True: Learn the Ultimate Language and Become a Better Programmer*.
No Starch Press, San Francisco, CA, USA, 1st edition, 2015.
URL: https://www.braveclojure.com/clojure-for-the-brave-and-true/.

📖 Simon Tatham.
Metaprogramming custom control structures in C.
https://www.chiark.greenend.org.uk/~sgtatham/mp/, 2012.
[Online; accessed 07-Feb-2018].