

Having Your Cake and Eating it Too: Combining Strong and Eventual Consistency

Konrad Siek and Paweł T. Wojciechowski

Poznań University of Technology

{konrad.siek,pawel.t.wojciechowski}@cs.put.edu.pl

13 IV 2014



Distributed Systems Group

<http://dsg.cs.put.poznan.pl>

Software Transactional Memory

```
def thread:  
    lock_a.acquire()  
    lock_b.acquire()  
    a = b  
    lock_a.release()  
    b = b + 1  
    lock_b.release()
```



Software Transactional Memory

```
def thread:  
    lock_a.acquire()  
    lock_b.acquire()  
    a = b  
    lock_a.release()  
    b = b + 1  
    lock_b.release()
```

```
def thread:  
    transaction.start()  
    a = b  
    b = b + 1  
    transaction.commit()
```

Software Transactional Memory

```
def thread:  
    lock_a.acquire()  
    lock_b.acquire()  
    a = b  
    lock_a.release()  
    b = b + 1  
    lock_b.release()
```

```
def thread:  
    transaction.start()  
    a = b  
    b = b + 1  
    transaction.commit()
```

Advantages:

- ease of use on top
- efficient concurrency control under the hood

Transaction Abstraction

Transaction:

$$T_i \llbracket op_1, op_2, \dots, op_n \rrbracket$$

where $op = \{ r(x)v, w(x)v, \dots \}$

and x is some shared object

Commitment:

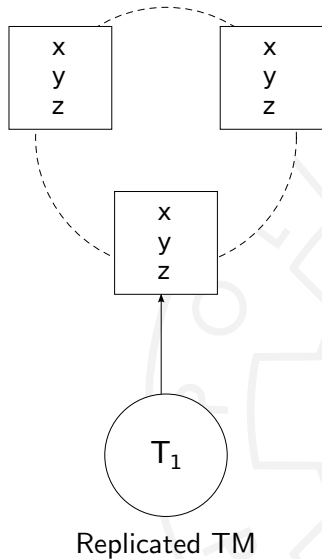
$$\{x = 1\} \quad T_i \llbracket w(x)2 \rrbracket \quad \{x = 2\}$$

Rollback:

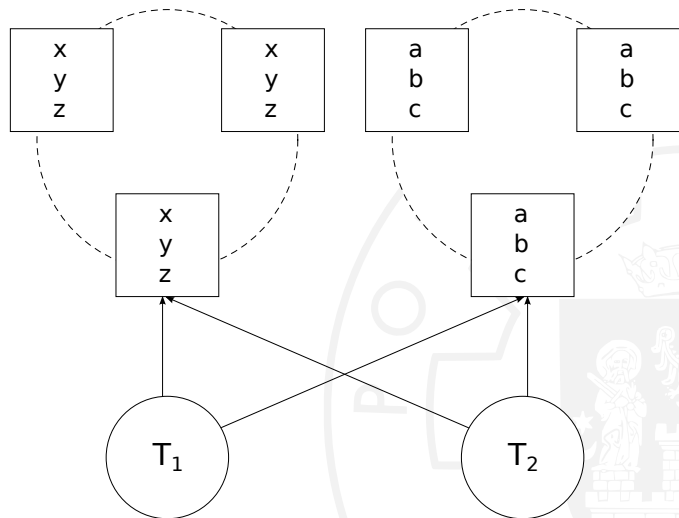
$$\{x = 1\} \quad T_i \llbracket w(x)2, \curvearrowright \rrbracket \quad \{x = 1\}$$

$$\{x = 1\} \quad T_i \llbracket w(x)2, \curvearrowright \rrbracket \rightarrow T'_i \llbracket w(x)2 \rrbracket \quad \{x = 2\}$$

Distributed Transactional Memory

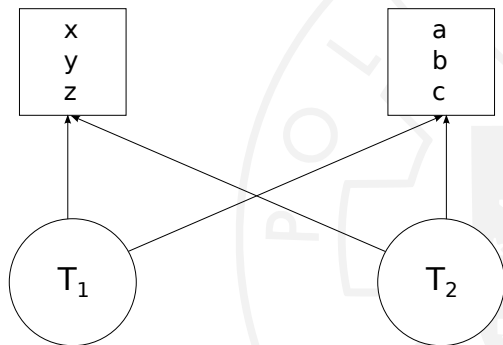


Distributed Transactional Memory



Distributed Transactions

Distributed Transactional Memory



Distributed Transactions

Supremum Versioning Algorithm

Pessimistic approach

$$\begin{array}{l} \{x = 1, y = 1\} \quad T_1 \llbracket r(x)1, w(x)2 \rrbracket \\ | T_2 \llbracket \phantom{\{x = 3, y = 2\}} \rrbracket \end{array} \quad \rightarrow r(x)2, w(x)3 \quad \{x = 3, y = 2\}$$

Supremum Versioning Algorithm

Pessimistic approach

$$\begin{array}{l} \{x = 1, y = 1\} \quad T_1 \llbracket r(x)1, w(x)2 \rrbracket \\ | T_2 \llbracket , w(x)3 \rrbracket \quad \{x = 3, y = 2\} \end{array}$$

(Note: In the original image, the $r(x)$ in the second transaction is highlighted in red.)

- Defer execution to prevent conflicts

Supremum Versioning Algorithm

Pessimistic approach

$$\begin{array}{l} \{x = 1, y = 1\} \quad T_1 \llbracket r(x)1, w(x)2 \rrbracket \\ | T_2 \llbracket \phantom{\{x = 3, y = 2\}} \rrbracket \end{array} \quad \{x = 3, y = 2\}$$

(Note: In the original image, the expression $r(x)2, w(x)3$ is highlighted in red and has an arrow pointing to it from the right.)

- Defer execution to prevent conflicts (tolerate high contention)

Supremum Versioning Algorithm

Pessimistic approach

$$\begin{array}{l} \{x = 1, y = 1\} \quad T_1 \llbracket r(x)1, w(x)2 \rrbracket \\ | T_2 \llbracket , w(x)3 \rrbracket \quad \{x = 3, y = 2\} \end{array}$$

- Defer execution to prevent conflicts (tolerate high contention)
- Avoid (most) forced aborts

Supremum Versioning Algorithm

Pessimistic approach

$$\begin{array}{l} \{x = 1, y = 1\} \quad T_1 \llbracket r(x)1, w(x)2 \rrbracket \\ | T_2 \llbracket \rightarrow r(x)2, w(x)3 \rrbracket \quad \{x = 3, y = 2\} \end{array}$$

- Defer execution to prevent conflicts (tolerate high contention)
- Avoid (most) forced aborts (safe irrevocable operations)

Early release on last use

$$\begin{array}{l} \{x = 1, y = 1\} \quad T_1 \llbracket r(x)1, w(x)2, r(y)1, w(y)2 \rrbracket \\ | T_2 \llbracket \rightarrow r(x)2, w(x)3 \rrbracket \quad \{x = 3, y = 2\} \end{array}$$

Supremum Versioning Algorithm

Pessimistic approach

$$\begin{array}{l} \{x = 1, y = 1\} \quad T_1 \llbracket r(x)1, w(x)2 \rrbracket \\ | T_2 \llbracket \rightarrow r(x)2, w(x)3 \rrbracket \quad \{x = 3, y = 2\} \end{array}$$

- Defer execution to prevent conflicts (tolerate high contention)
- Avoid (most) forced aborts (safe irrevocable operations)

Early release on last use

$$\begin{array}{l} \{x = 1, y = 1\} \quad T_1 \llbracket r(x)1, w(x)2, r(y)1, w(y)2 \rrbracket \\ | T_2 \llbracket \rightarrow r(x)2, w(x)3 \rrbracket \quad \{x = 3, y = 2\} \end{array}$$

Bank Application



Bank Application

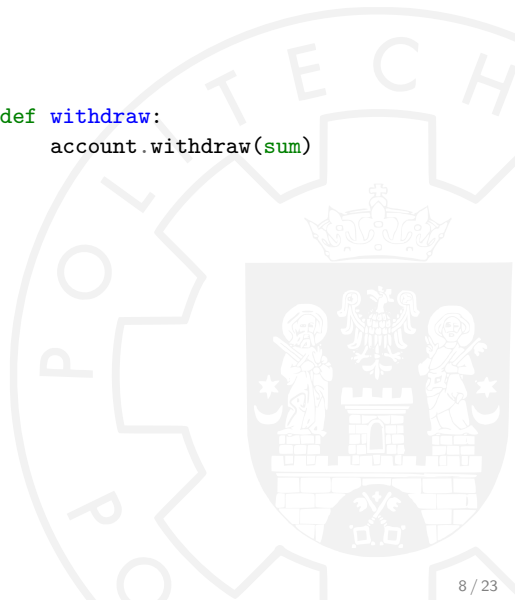
```
def deposit:  
    account.deposit(sum)
```



Bank Application

```
def deposit:  
    account.deposit(sum)
```

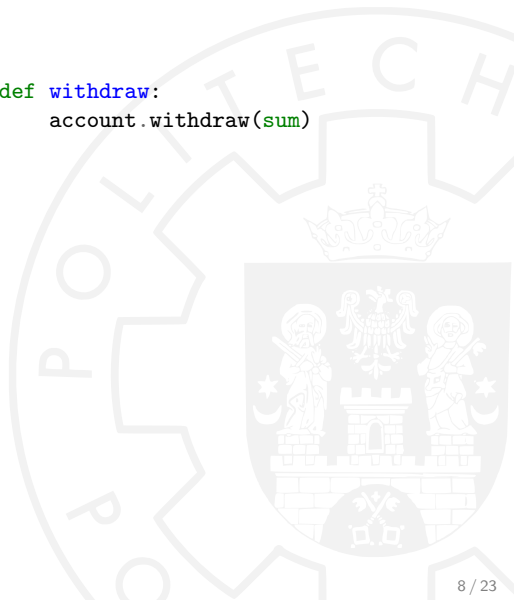
```
def withdraw:  
    account.withdraw(sum)
```



Bank Application

```
def deposit:  
    account.deposit(sum)  
  
def balance:  
    print account.getBalance()
```

```
def withdraw:  
    account.withdraw(sum)
```



Bank Application

```
def deposit:  
    account.deposit(sum)
```

```
def balance:  
    print account.getBalance()
```

```
def withdraw:  
    account.withdraw(sum)
```

```
def transfer:  
    account1.withdraw(sum)  
    account2.deposit(sum)
```

Bank Application

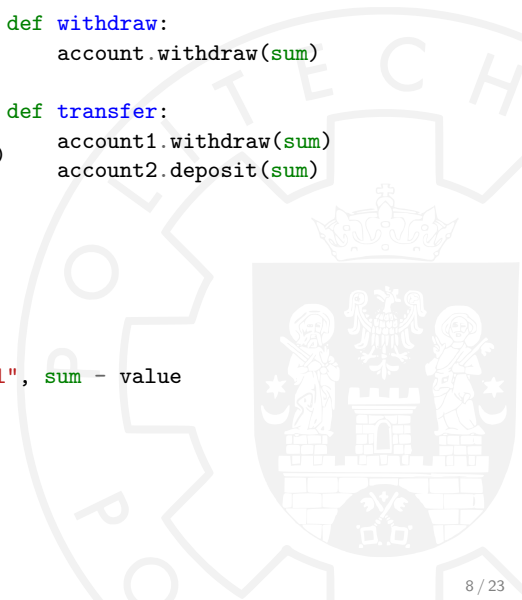
```
def deposit:  
    account.deposit(sum)
```

```
def balance:  
    print account.getBalance()
```

```
def audit:  
    for a in accounts:  
        sum += a.getBalance()  
    value = bank.getCapital()  
    bank.setCapital(sum)  
    print "Accumulated capital", sum - value
```

```
def withdraw:  
    account.withdraw(sum)
```

```
def transfer:  
    account1.withdraw(sum)  
    account2.deposit(sum)
```



Bank Application

```
def deposit:  
    transaction.start()  
    account.deposit(sum)  
    transaction.commit()
```

```
def balance:  
    transaction.start()  
    print account.getBalance()  
    transaction.commit()
```

```
def audit:  
    transaction.start()  
    for a in accounts:  
        sum += a.getBalance()  
    value = bank.getCapital()  
    bank.setCapital(sum)  
    print "Accumulated capital", sum - value  
    transaction.commit()
```

```
def withdraw:  
    transaction.start()  
    account.withdraw(sum)  
    transaction.commit()
```

```
def transfer:  
    transaction.start()  
    account1.withdraw(sum)  
    account2.deposit(sum)  
    transaction.commit()
```


Bank Application

```
def deposit:  
    transaction.start()  
    account.deposit(sum)  
    transaction.commit()
```

```
def balance:  
    transaction.start()  
    print account.getBalance()  
    transaction.commit()
```

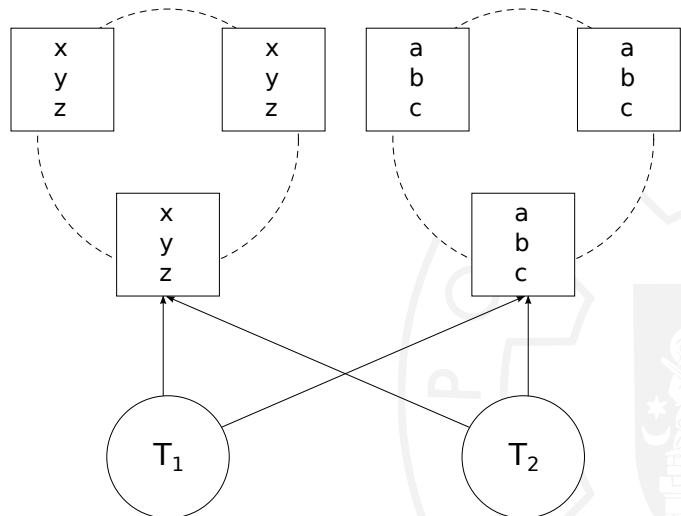
```
def audit:  
    transaction.start()  
    for a in accounts:  
        sum += a.getBalance()  
    value = bank.getCapital()  
    bank.setCapital(sum)  
    print "Accumulated capital", sum - value  
    transaction.commit()
```

```
def withdraw:  
    transaction.start()  
    account.withdraw(sum)  
    transaction.commit()
```

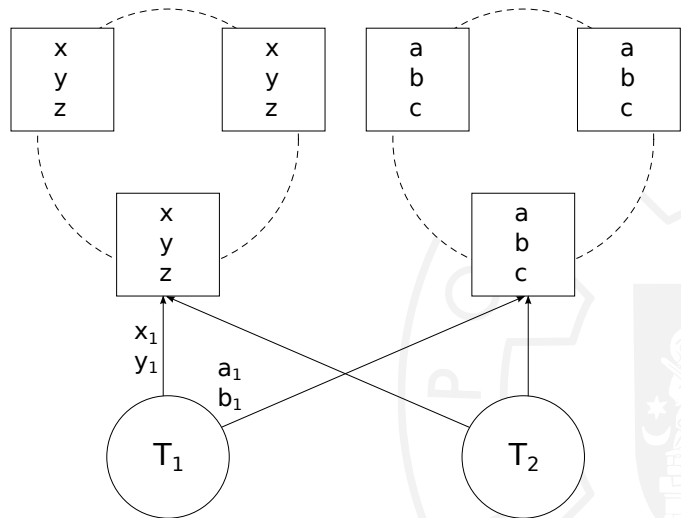
```
def transfer:  
    transaction.start()  
    account1.withdraw(sum)  
    account2.deposit(sum)  
    transaction.commit()
```

weaken consistency → improve efficiency

Eventually Consistent Extension



Eventually Consistent Extension



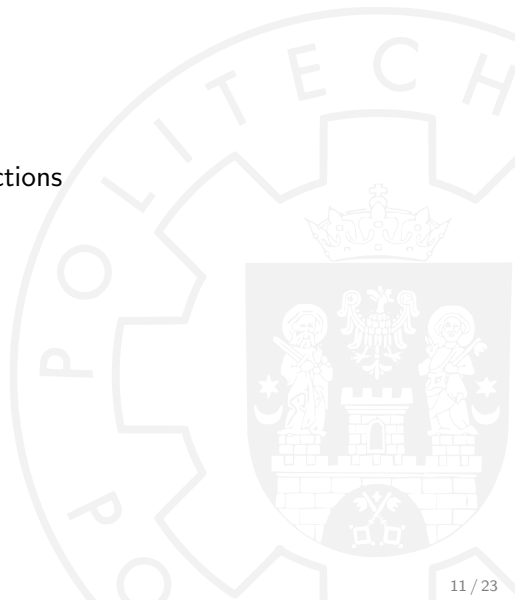
What we require from weak transactions

- do not wait for variables



What we require from weak transactions

- do not wait for variables
- do not block other transactions



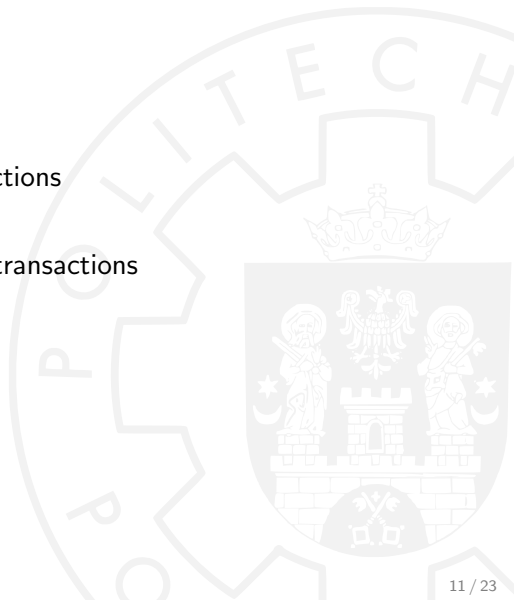
What we require from weak transactions

- do not wait for variables
- do not block other transactions
- internal consistency



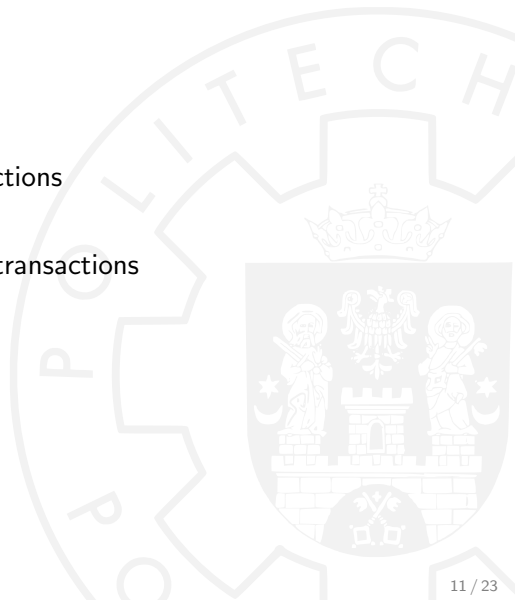
What we require from weak transactions

- do not wait for variables
- do not block other transactions
- internal consistency
- do not disturb consistent transactions



What we require from weak transactions

- do not wait for variables
- do not block other transactions
- internal consistency
- do not disturb consistent transactions
- converge



Transaction Modes

Transaction T_1

$$T_1 \llbracket r(x)v_c, w(x)u_c \rrbracket$$



Transaction Modes

Transaction T_1

$$T_1 \llbracket r(x)v_c, w(x)u_c \rrbracket$$

Consistent mode

$$T_1^c \llbracket r(x)v_c, w(x)u_c \rrbracket$$



Transaction Modes

Transaction T_1

$$T_1 \llbracket r(x)v_c, w(x)u_c \rrbracket$$

Consistent mode

$$T_1^c \llbracket r(x)v_c, w(x)u_c \rrbracket$$

Eventually consistent mode

$$T_1^{ec} \llbracket r(x)v_{ec}, w(x)u_{ec} \rrbracket$$



Transaction Modes

Transaction T_1

$$T_1 \llbracket r(x)v_c, w(x)u_c \rrbracket$$

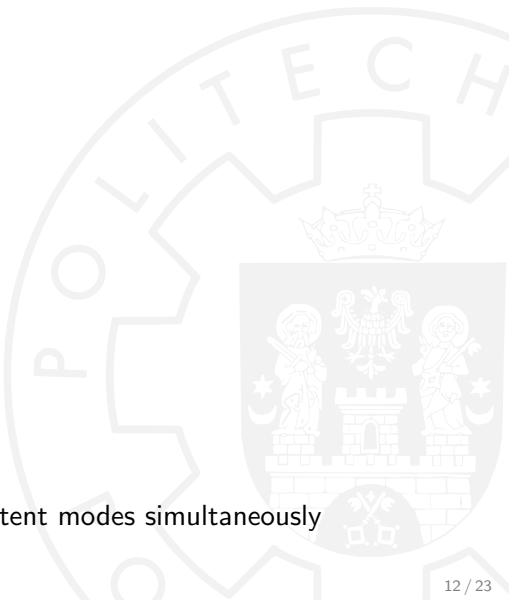
Consistent mode

$$T_1^c \llbracket r(x)v_c, w(x)u_c \rrbracket$$

Eventually consistent mode

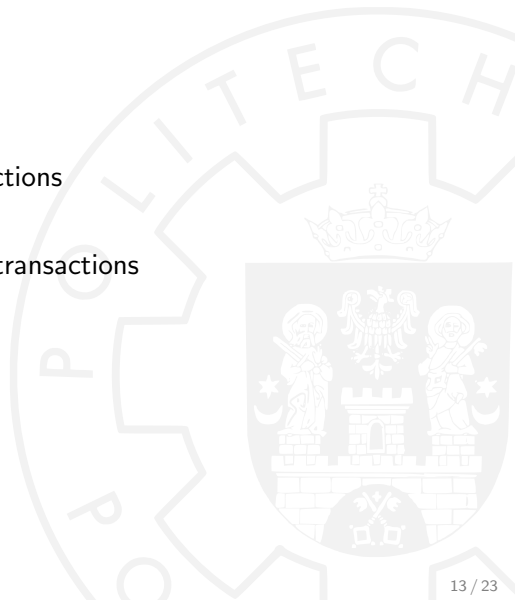
$$T_1^{ec} \llbracket r(x)v_{ec}, w(x)u_{ec} \rrbracket$$

Execute consistent and inconsistent modes simultaneously



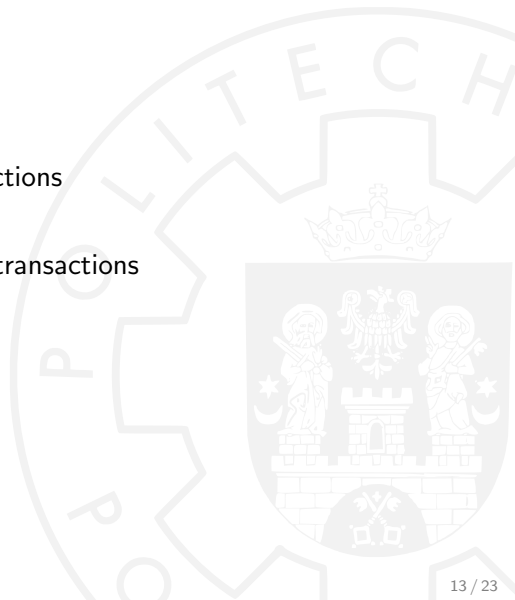
What we require from weak transactions

- do not wait for variables
- do not block other transactions
- internal consistency
- do not disturb consistent transactions
- converge



What we require from weak transactions

- do not wait for variables
- do not block other transactions
- internal consistency
- do not disturb consistent transactions
- converge ✓

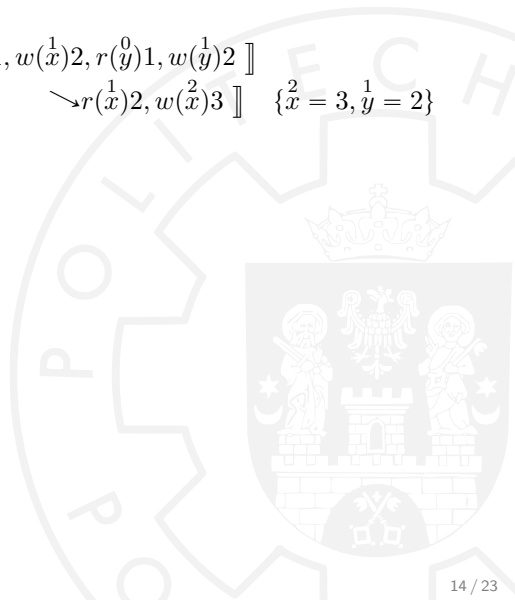


Variable Modification Versions

$$\{x = 1, y = 1\} \quad T_1 \left[\begin{array}{l} r(x)1, w(x)2, r(y)1, w(y)2 \\ | T_2 \left[\begin{array}{l} \rightarrow r(x)2, w(x)3 \end{array} \right] \end{array} \right] \quad \{x = 3, y = 2\}$$

Variable Modification Versions

$$\{x^0 = 1, y^0 = 1\} \quad T_1^{ec} \left[\left[r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \right] \right. \\ \left. \mid T_2 \left[\begin{array}{l} \rightarrow r(x^1)2, w(x^2)3 \end{array} \right] \quad \{x^2 = 3, y^1 = 2\} \right]$$



Variable Modification Versions

$$\{x^0 = 1, y^0 = 1\} \quad T_1^{ec} \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \quad \{x^2 = 3, y^1 = 2\}$$

Snapshot Read Consistency

$$T_1 \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2, w(y^2)3 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \\ | T_3 \llbracket \searrow r(x^2)3, w(x^3)4, r(y^2)3, w(y^3)4 \rrbracket \\ | T_4 \llbracket \searrow r(x^3)4, w(x^4)5 \rrbracket$$

Variable Modification Versions

$$\{x^0 = 1, y^0 = 1\} \quad T_1^{ec} \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \quad \{x^2 = 3, y^1 = 2\}$$

Snapshot Read Consistency

$$T_1 \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2, w(y^2)3 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \\ | T_3 \llbracket \searrow r(x^2)3, w(x^3)4, r(y^2)3, w(y^3)4 \rrbracket \\ | T_4 \llbracket \searrow r(x^3)4, w(x^4)5 \rrbracket \\ \{x^1, y^2\}$$

Variable Modification Versions

$$\{x^0 = 1, y^0 = 1\} \quad T_1^{ec} \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \quad \{x^2 = 3, y^1 = 2\}$$

Snapshot Read Consistency

$$T_1 \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2, w(y^2)3 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \\ | T_3 \llbracket \searrow r(x^2)3, w(x^3)4, r(y^2)3, w(y^3)4 \rrbracket \\ | T_4 \llbracket \searrow r(x^3)4, w(x^4)5 \rrbracket \\ \{x^1, y^2\}, \{x^2, y^2\}$$

Variable Modification Versions

$$\{x^0 = 1, y^0 = 1\} \quad T_1^{ec} \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \quad \{x^2 = 3, y^1 = 2\}$$

Snapshot Read Consistency

$$T_1 \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2, w(y^2)3 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \\ | T_3 \llbracket \searrow r(x^2)3, w(x^3)4, r(y^2)3, w(y^3)4 \rrbracket \\ | T_4 \llbracket \searrow r(x^3)4, w(x^4)5 \rrbracket$$

$$\{x^1, y^2\}, \{x^2, y^2\}, \{x^3, y^3\}, \{x^4, y^3\}$$

Variable Modification Versions

$$\{x^0 = 1, y^0 = 1\} \quad T_1^{ec} \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \quad \{x^2 = 3, y^1 = 2\}$$

Snapshot Read Consistency

$$T_1 \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2, w(y^2)3 \rrbracket \\ | T_2 \llbracket \searrow r(x^1)2, w(x^2)3 \rrbracket \\ | T_3 \llbracket \searrow r(x^2)3, w(x^3)4, r(y^2)3, w(y^3)4 \rrbracket \\ | T_4 \llbracket \searrow r(x^3)4, w(x^4)5 \rrbracket$$

$$\{x^1, y^2\}, \{x^2, y^2\}, \{x^3, y^3\}, \{x^4, y^3\}, \{x^3, y^2\}$$

Variable Modification Versions

$$\{x^0 = 1, y^0 = 1\} \quad T_1^{ec} \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \rrbracket \\ | T_2 \llbracket , , r(x^1)2, w(x^2)3 \rrbracket \quad \{x^2 = 3, y^1 = 2\}$$

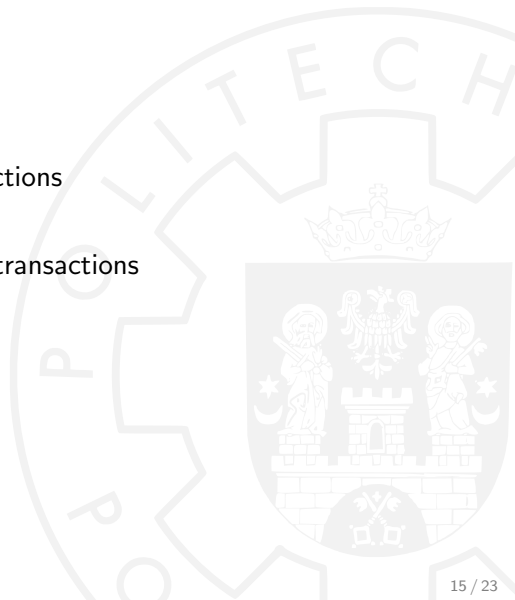
Snapshot Read Consistency

$$T_1 \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2, w(y^2)3 \rrbracket \\ | T_2 \llbracket , , r(x^1)2, w(x^2)3 \rrbracket \\ | T_3 \llbracket , , , w(x^2)3, w(x^3)4, r(y^2)3, w(y^3)4 \rrbracket \\ | T_4 \llbracket , , , , w(x^3)4, w(x^4)5 \rrbracket$$

$$\{x^1, y^2\}, \{x^2, y^2\}, \{x^3, y^3\}, \{x^4, y^3\}, \{x^3, y^2\}, \{x^*, y^1\}$$

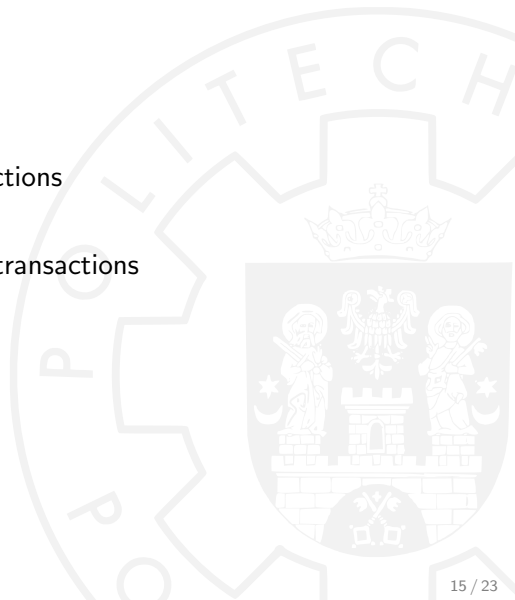
What we require from weak transactions

- do not wait for variables
- do not block other transactions
- internal consistency
- do not disturb consistent transactions
- converge ✓



What we require from weak transactions

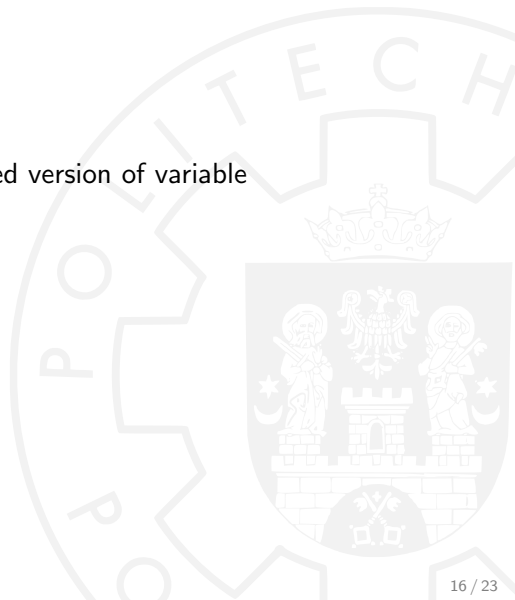
- do not wait for variables
- do not block other transactions
- internal consistency ✓
- do not disturb consistent transactions
- converge ✓



Consistent Snapshot in Practice

Transactions:

- record the latest committed version of variable



Consistent Snapshot in Practice

Transactions:

- record the latest committed version of variable
- record the latest released version of variable (early release)

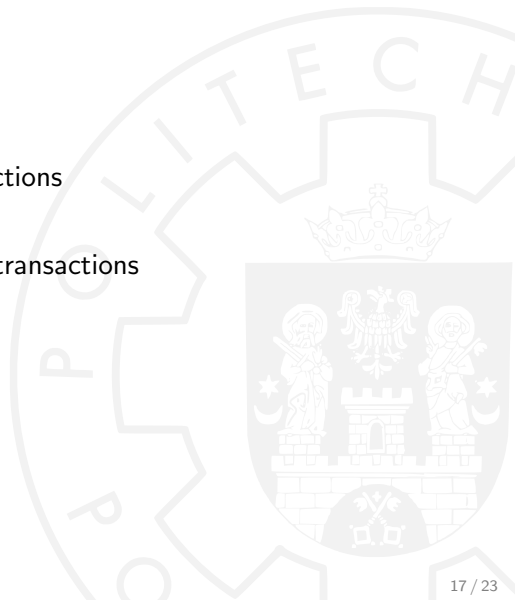
Consistent Snapshot in Practice

Transactions:

- record the latest committed version of variable
- record the latest released version of variable (early release)
- when releasing a variable early: record variables that were not released early

What we require from weak transactions

- do not wait for variables
- do not block other transactions
- internal consistency ✓
- do not disturb consistent transactions
- converge ✓

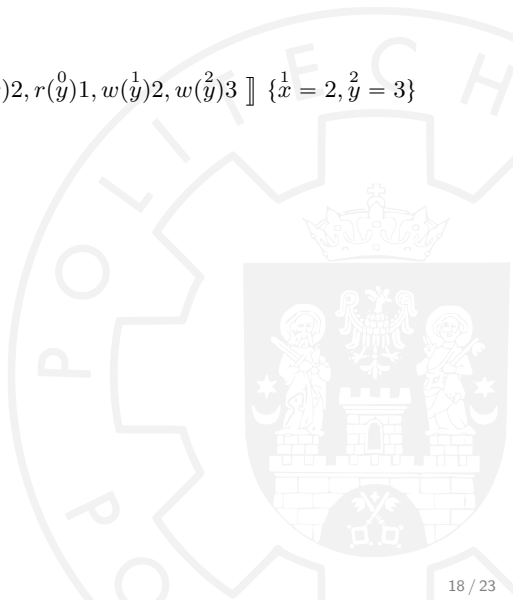


What we require from weak transactions

- do not wait for variables ✓
- do not block other transactions ✓
- internal consistency ✓
- do not disturb consistent transactions
- converge ✓

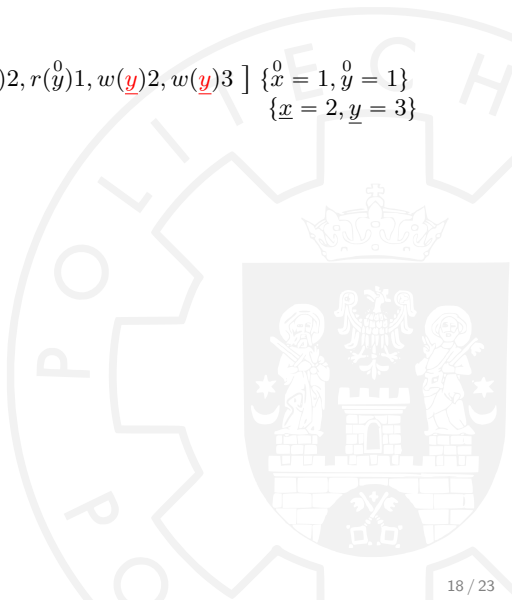
Write Bufferring

$$\{x^0 = 1, y^0 = 1\} T_1 \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2, w(y^2)3 \rrbracket \{x^1 = 2, y^2 = 3\}$$



Write Bufferring

$$\{x^0 = 1, y^0 = 1\} T_1 [r(\underline{x})1, w(\underline{x})2, r(\underline{y})1, w(\underline{y})2, w(\underline{y})3] \{x^0 = 1, y^0 = 1\}$$
$$\{x = 2, y = 3\}$$



Write Buffering

$$\{x^0 = 1, y^0 = 1\} T_1 [r(\underline{x})1, w(\underline{x})2, r(\underline{y})1, w(\underline{y})2, w(\underline{y})3] \{x^0 = 1, y^0 = 1\} \\ \{x = 2, y = 3\}$$

Consistent mode either:

- applies the buffered writes

Write Buffering

$$\{x^0 = 1, y^0 = 1\} T_1 [r(\underline{x})1, w(\underline{x})2, r(\underline{y})1, w(\underline{y})2, w(\underline{y})3] \{x^0 = 1, y^0 = 1\} \\ \{x = 2, y = 3\}$$

Consistent mode either:

- applies the buffered writes (if consistency condition allows)

Write Bufferring

$$\{x^0 = 1, y^0 = 1\} T_1 [r(\underline{x})1, w(\underline{x})2, r(\underline{y})1, w(\underline{y})2, w(\underline{y})3] \{x^0 = 1, y^0 = 1\}$$
$$\{x = 2, y = 3\}$$

Consistent mode either:

- applies the buffered writes (if consistency condition allows)
- re-executes from scratch

What we require from weak transactions

- do not wait for variables ✓
- do not block other transactions ✓
- internal consistency ✓
- do not disturb consistent transactions
- converge ✓

What we require from weak transactions

- do not wait for variables ✓
- do not block other transactions ✓
- internal consistency ✓
- do not disturb consistent transactions ✓
- converge ✓

Eventually Consistent SVA

$$\begin{aligned} \{x^0 = 1, y^0 = 1\} & \quad T_1 \left[\left[r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \right] \right. \\ & \quad | T_2^c \left[\begin{array}{l} \\ \searrow r(x^1)2, w(x^2)3 \end{array} \right] \quad \{x^2 = 3, y^1 = 2\} \\ & \quad | T_2^{ec} \left[r(x^0)1, w(\underline{x})2 \right] \end{aligned}$$

Eventually Consistent SVA

$$\begin{array}{l} \{x^0 = 1, y^0 = 1\} \quad T_1 \llbracket r(x^0)1, w(x^1)2, r(y^0)1, w(y^1)2 \rrbracket \\ | T_2^c \llbracket , w(x^2)3 \rrbracket \quad \{x^2 = 3, y^1 = 2\} \\ | T_2^{ec} \llbracket r(x^0)1, w(\underline{x})2 \rrbracket \\ | T_3 \llbracket r(x^2)3, w(x^3)4 \rrbracket \end{array}$$

Summary

- eventual consistency extension for pessimistic distributed TM
- minimal extra cost
- eventually consistent transactions read consistent snapshots
- strongly consistent transactions are unaffected

Related Papers:

Konrad Siek, Paweł T. Wojciechowski. *Brief Announcement: Towards a Fully-Articulated Pessimistic Distributed Transactional Memory*. In Proceedings of SPAA 2013: the 25th ACM Symposium on Parallelism in Algorithms and Architectures. July 2013.

Paweł T. Wojciechowski, Olivier Rütli and André Schiper. *SAMOA: A Framework for a Synchronisation-Augmented Microprotocol Approach*. In the Proceedings of IPDPS 2004: the 18th IEEE Parallel and Distributed Processing Symposium. April 2004.

?

