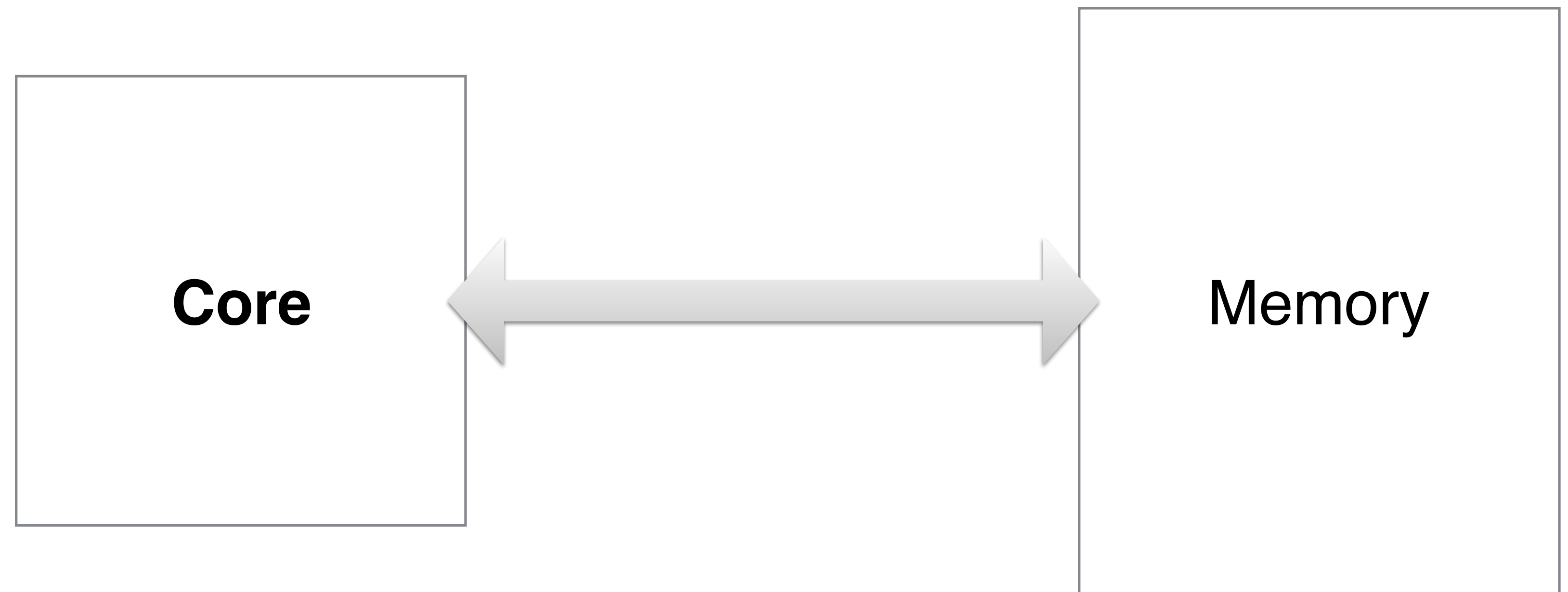# Real-time scheduling of STM transactions on multi-core platforms

**António Barros**, Patrick Meumeu Yomsi, Luís Miguel Pinho
CISTER seminar series
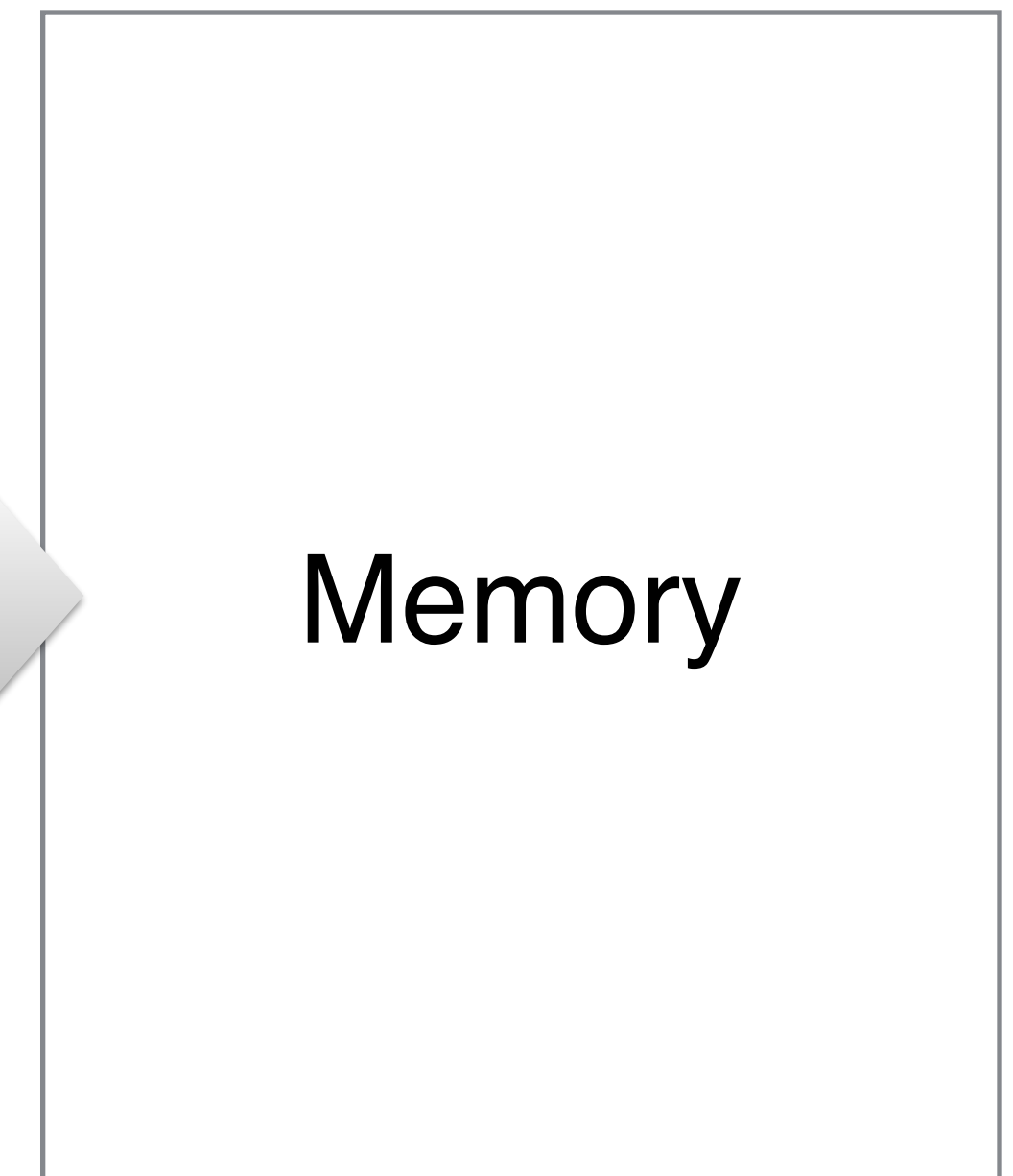xxth February 2015

# The problem

**Single core**

# The problem
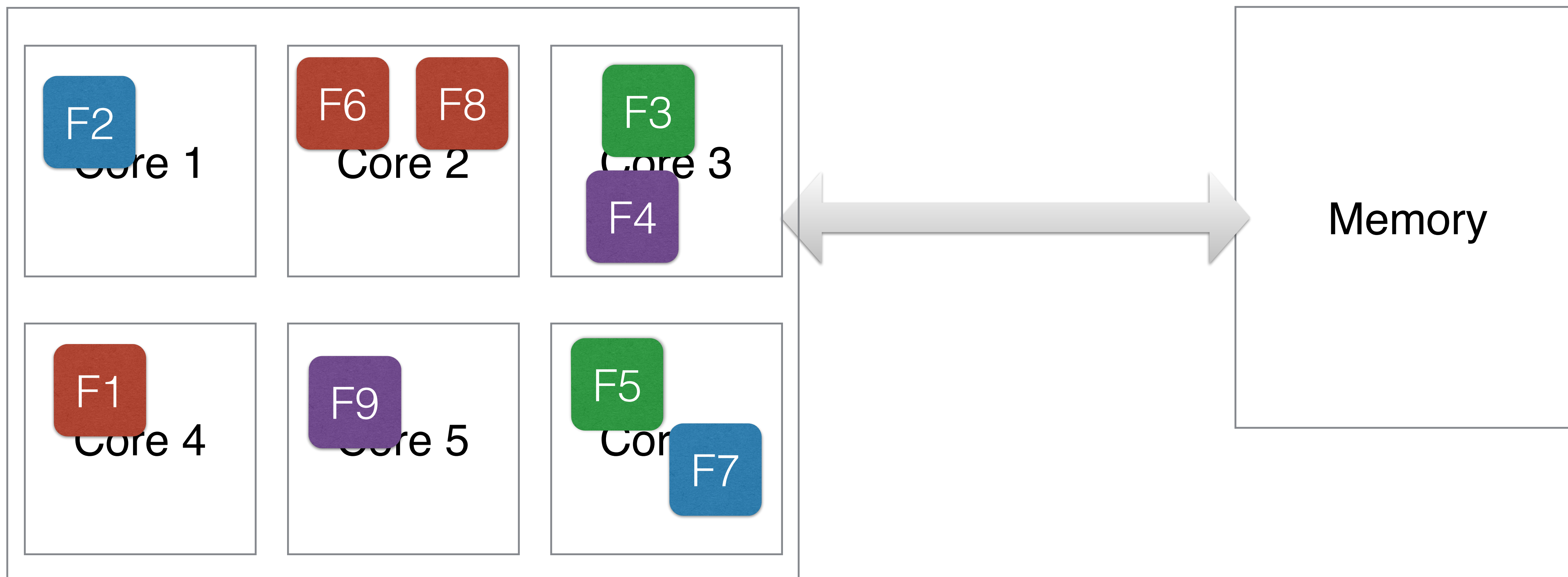
**Single core**



Memory

# The problem

Well understood theory and practice on unicore platforms!

Memory

# The problem

**Multi-core**

# The problem

Current and future embedded architectures...
- Multiple cores (tens, hundreds,...)
- No cache-coherency
- Single memory bus

F2

Core

F1

Core

F7

ory

# The problem

F2

Core

F1

Core

ory

Current and future embedded architectures...
- Multiple cores (tens, hundreds,...)
- No cache-coherency
- Single memory bus

Maybe OK (?) for sets of independent tasks...

# The problem

F2

Core

F1

Core

Current and future embedded architectures…
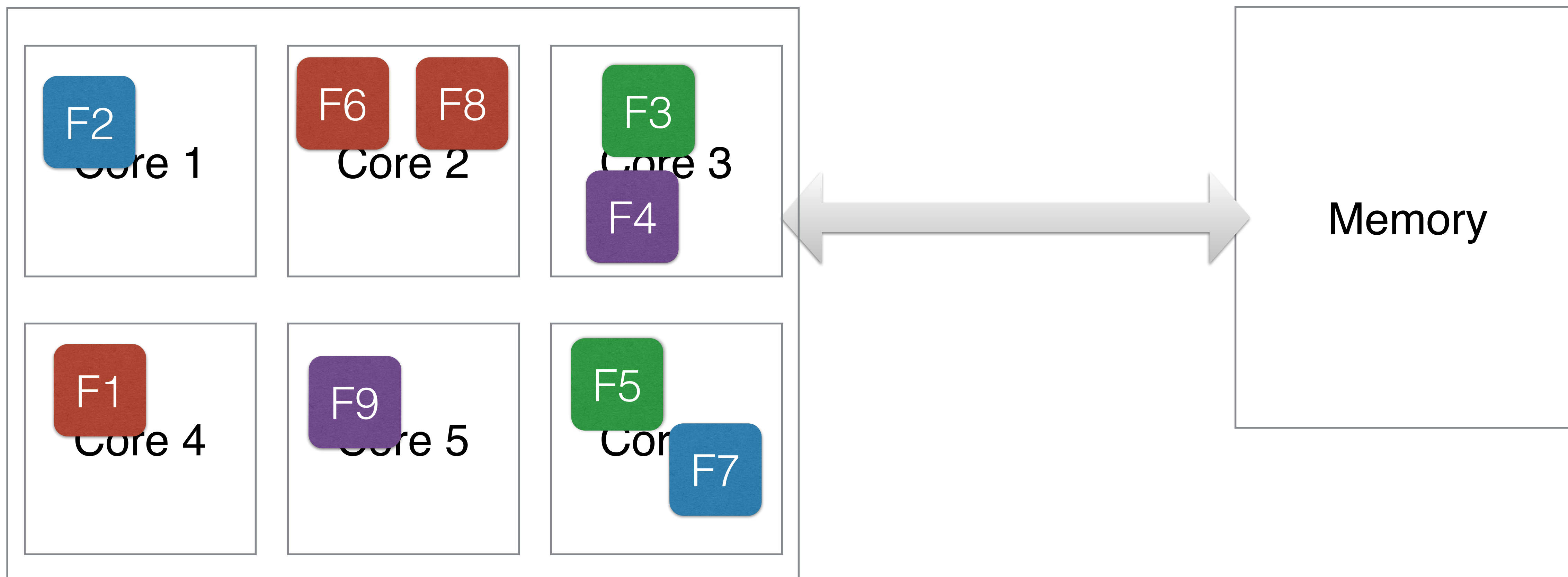- Multiple cores (tens, hundreds,…)
- No cache-coherency
- Single memory bus

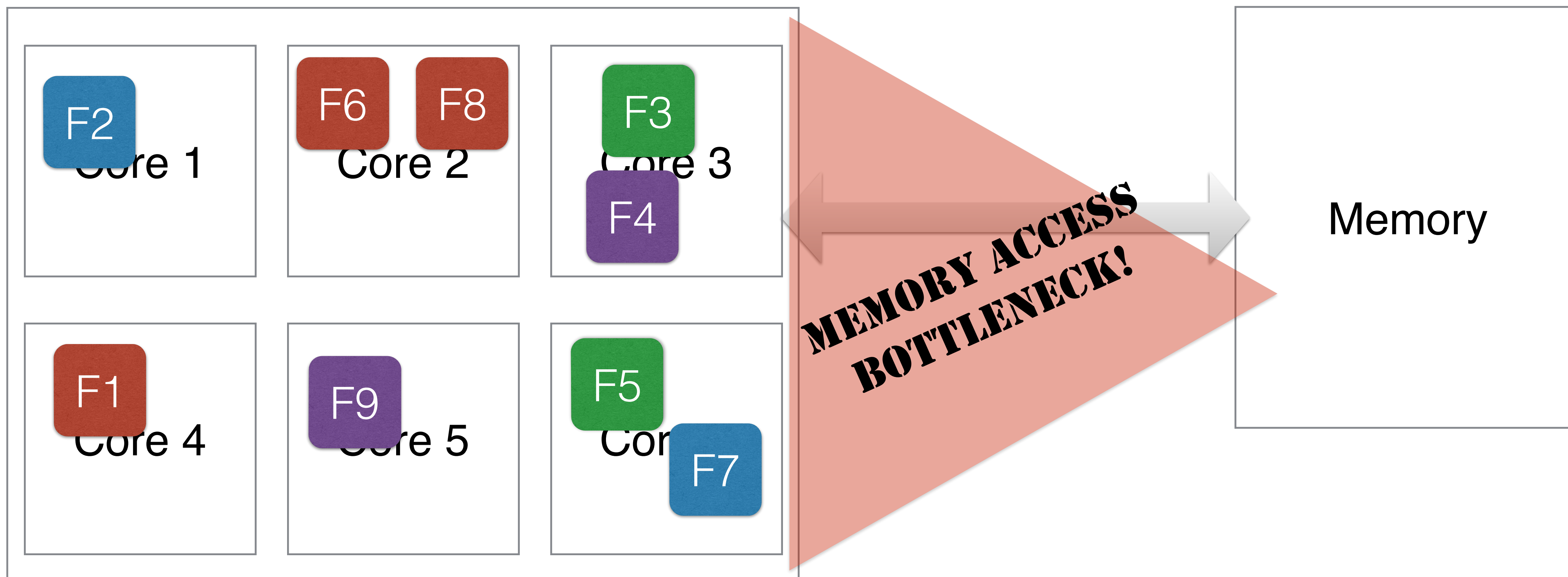Maybe OK (?) for sets of independent tasks…

What if tasks ARE NOT independent?
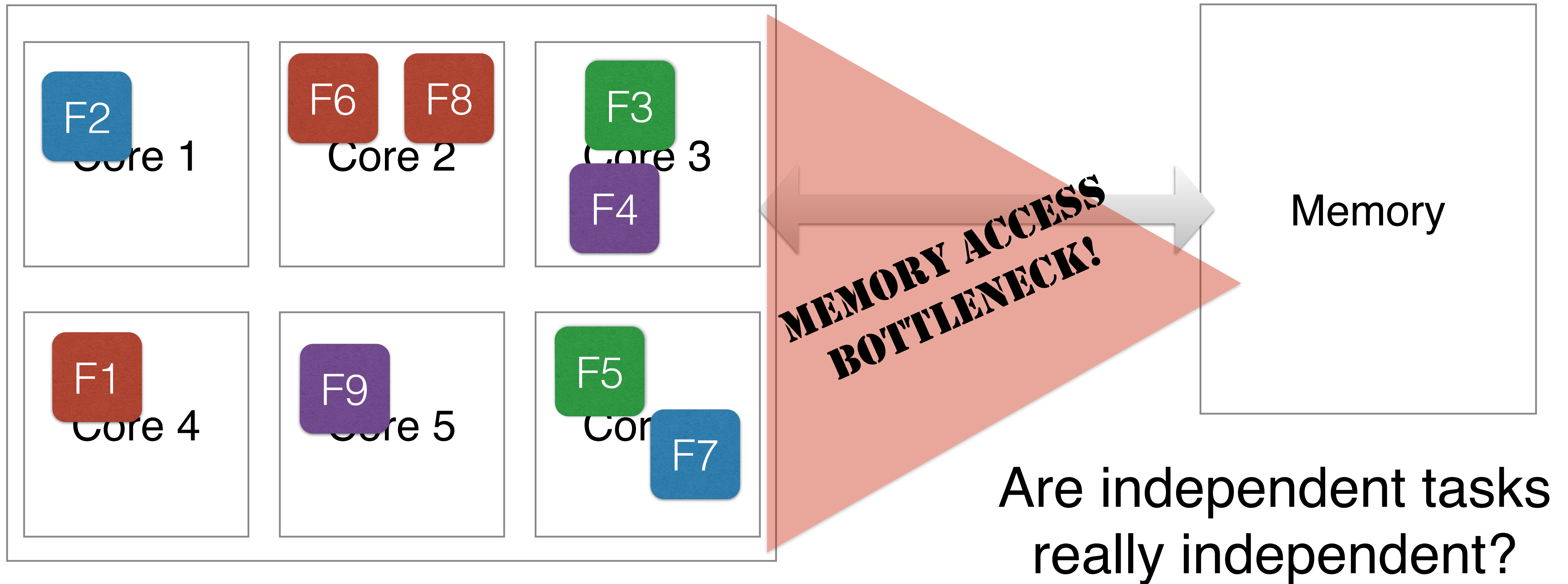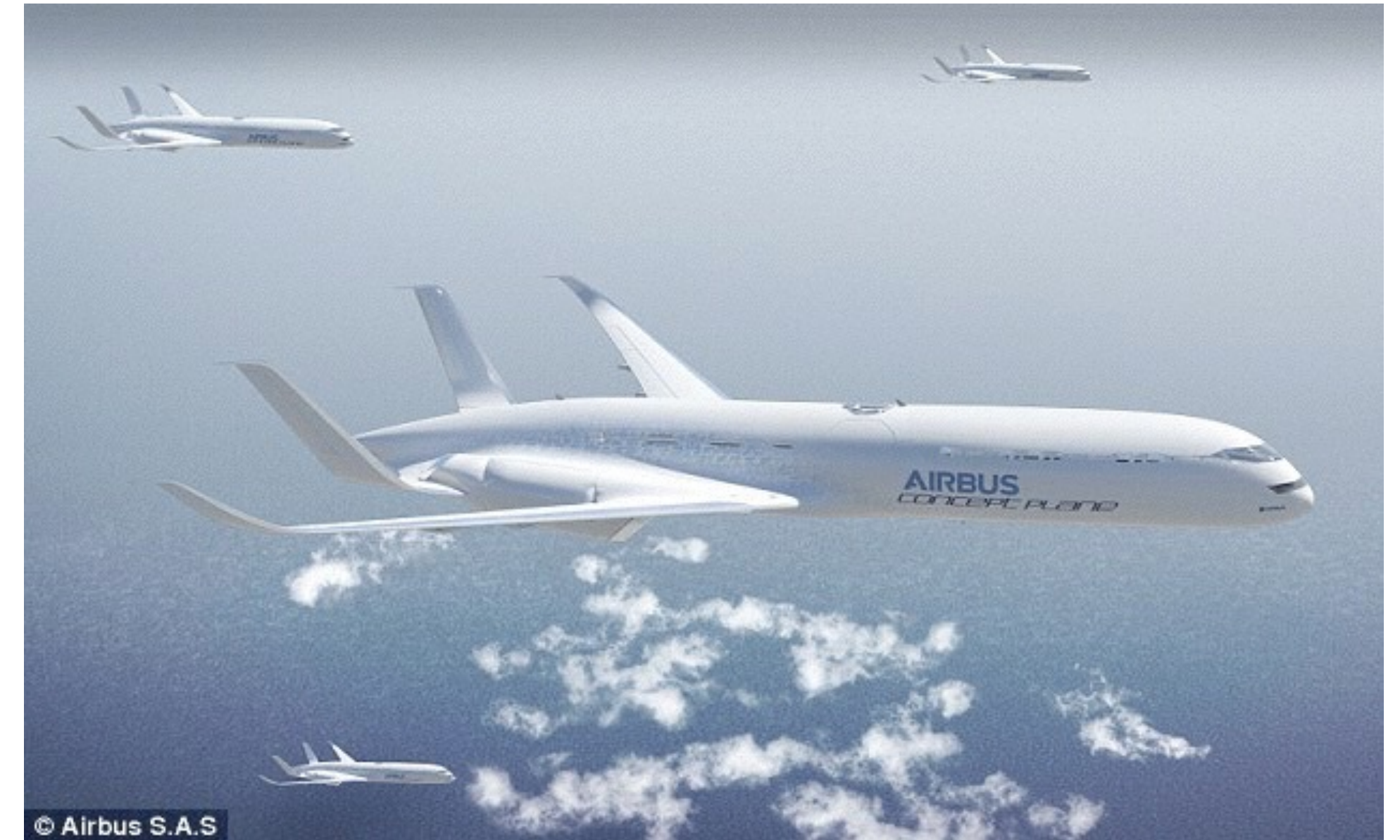
# The problem

**Multi-core**

# The problem

**Multi-core**

| | | |
|---|---|---|
| F2 — Core 1 | F6 F8 — Core 2 | F3 F4 — Core 3 |
| F1 — Core 4 | F9 — Core 5 | F5 F7 — Core 6 |

**MEMORY ACCESS BOTTLENECK!**

Memory

# The problem

**Multi-core**

Core 1
F2

Core 2
F6 F8

Core 3
F3
F4

Core 4
F1

Core 5
F9

Core 6
F5
F7

MEMORY ACCESS BOTTLENECK!

Memory

Are independent tasks really independent?

# Practical case: DO-178C

- New programming paradigm (enhancing explicit dependencies between functionalities).

- Spatial and temporal isolation among functionalities, depending on their criticality.

- Functionalities must be statically assigned to cores.

- Data dependencies must be mapped.



© Airbus S.A.S

# Attempted solutions

Recent proposal: FMLP*

- Global resources can be short or long (designer's choice), depending on length of critical sections

  - When blocked: busy waits on short, suspends on long

- Nested requests dictates joining resources into resource groups

  - one lock (queue lock or semaphore) per group

  - exclusively short- and long-groups

- Critical section code is executed non-preemptively

* A. Block, H. Leontyev, B. Brandenburg, and J. Anderson. A flexible real-time locking protocol for multiprocessors. In Proceedings of RTCSA 2007, pages 71–80, 2007.

# Attempted solutions

Our idea: STM + SRP-TM

- No groups and locks (at least, seen by the programmer)

- Contention is checked at run-time: just-in-time parallelism

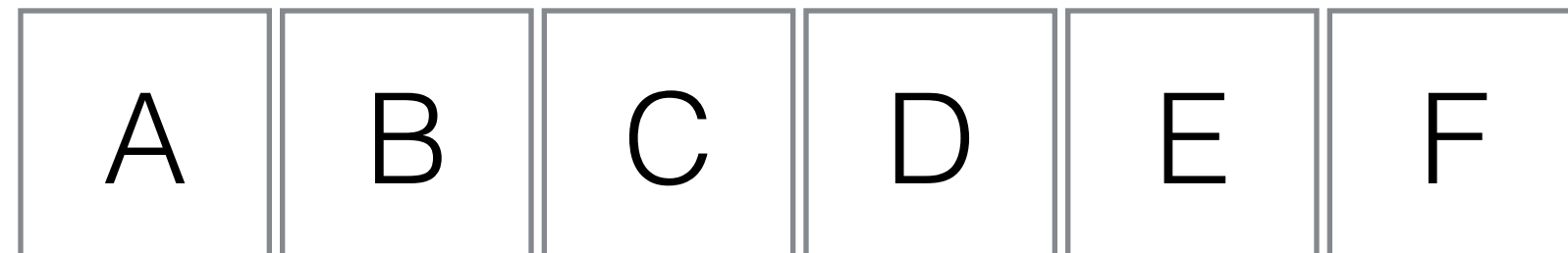- Upper-bound atomic section response time

- Limit blocking times

# Background
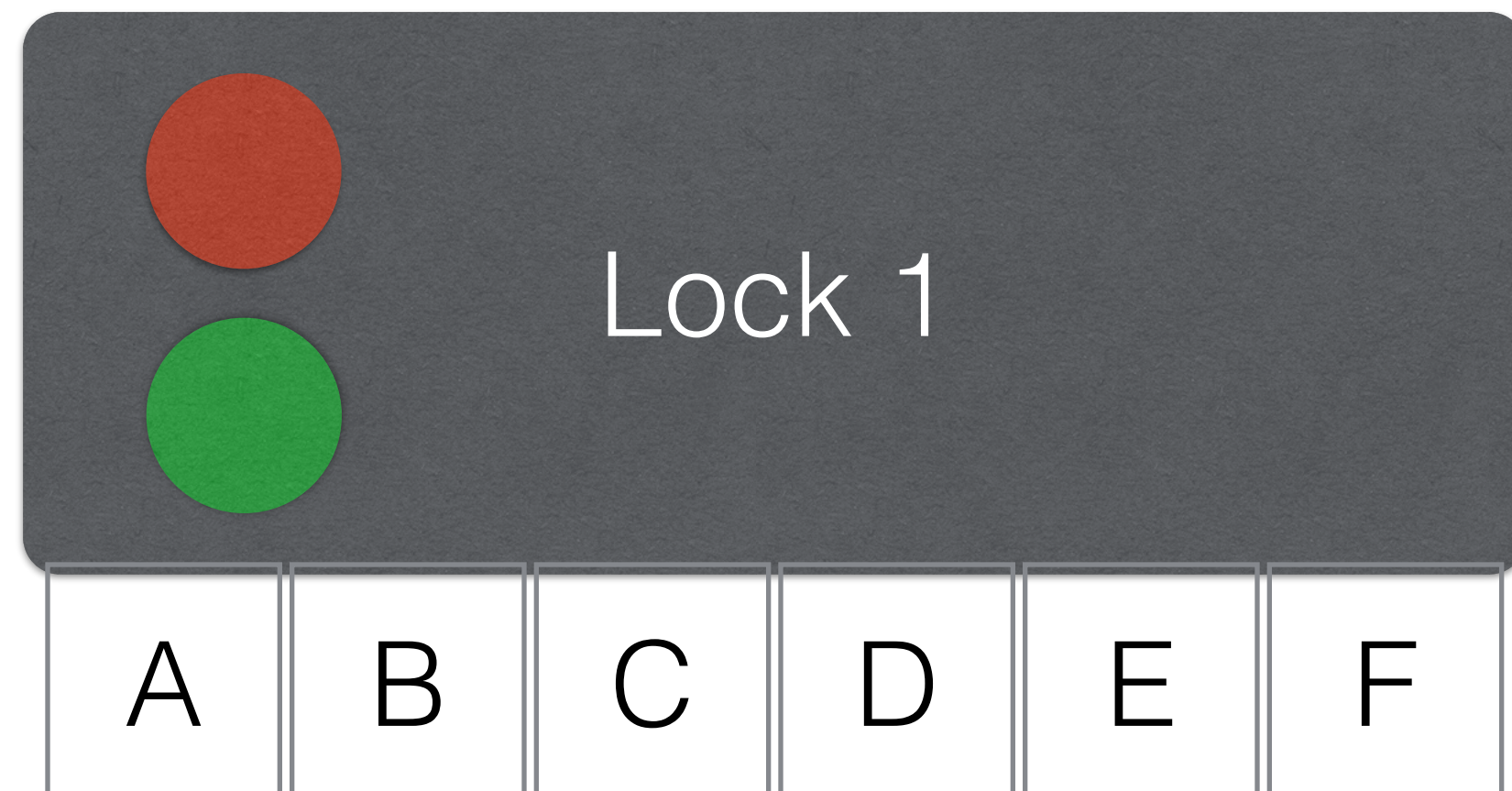
# Locks

- Coarse-grained locking
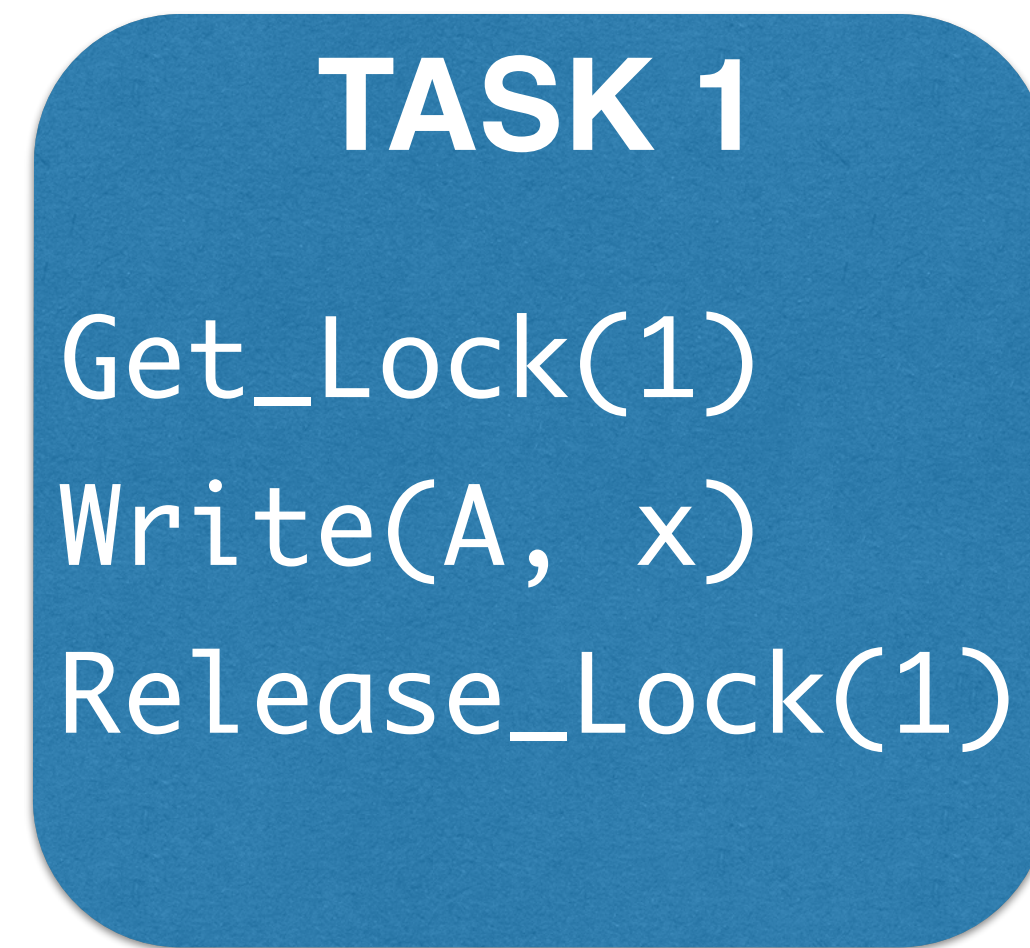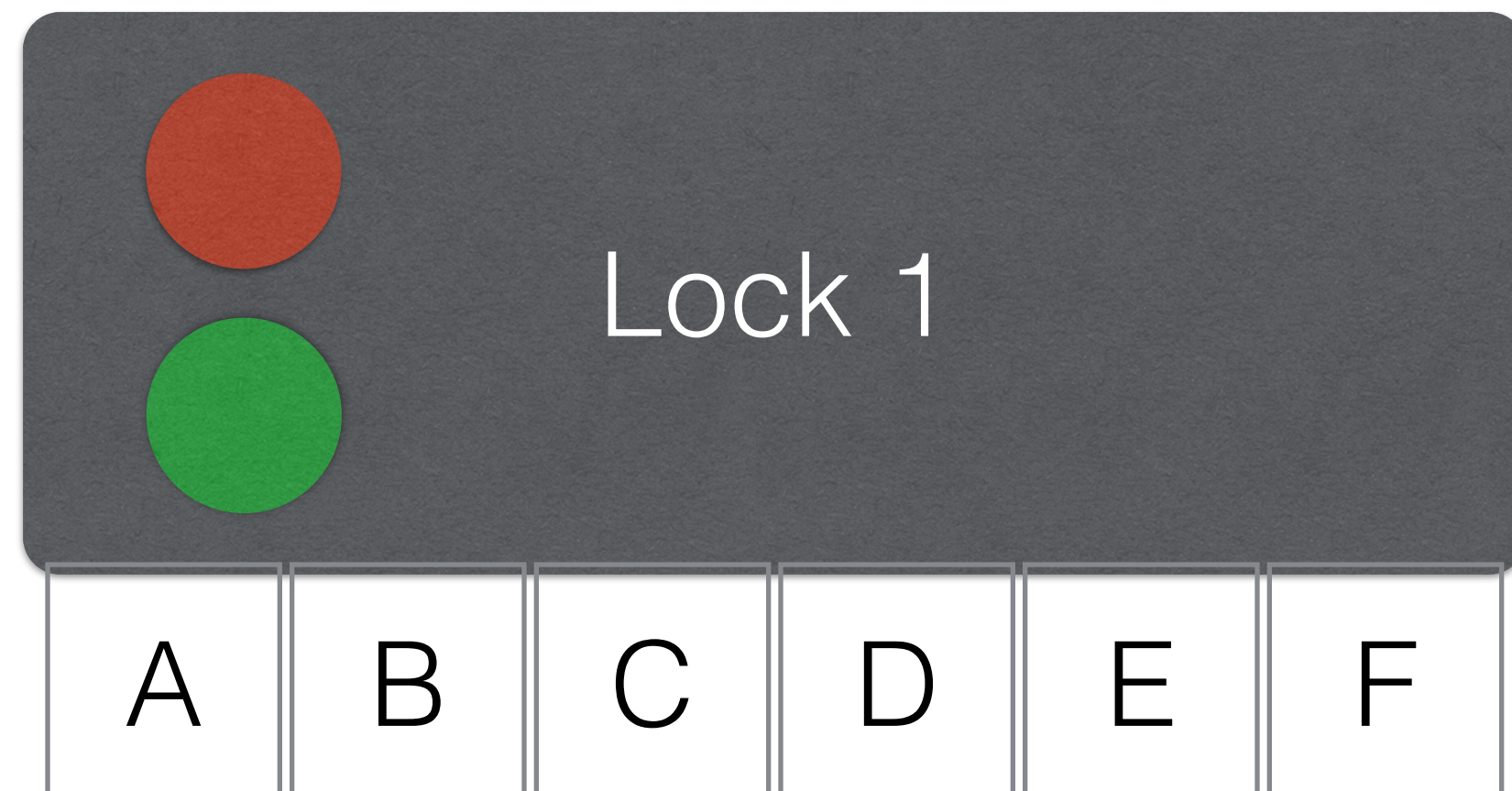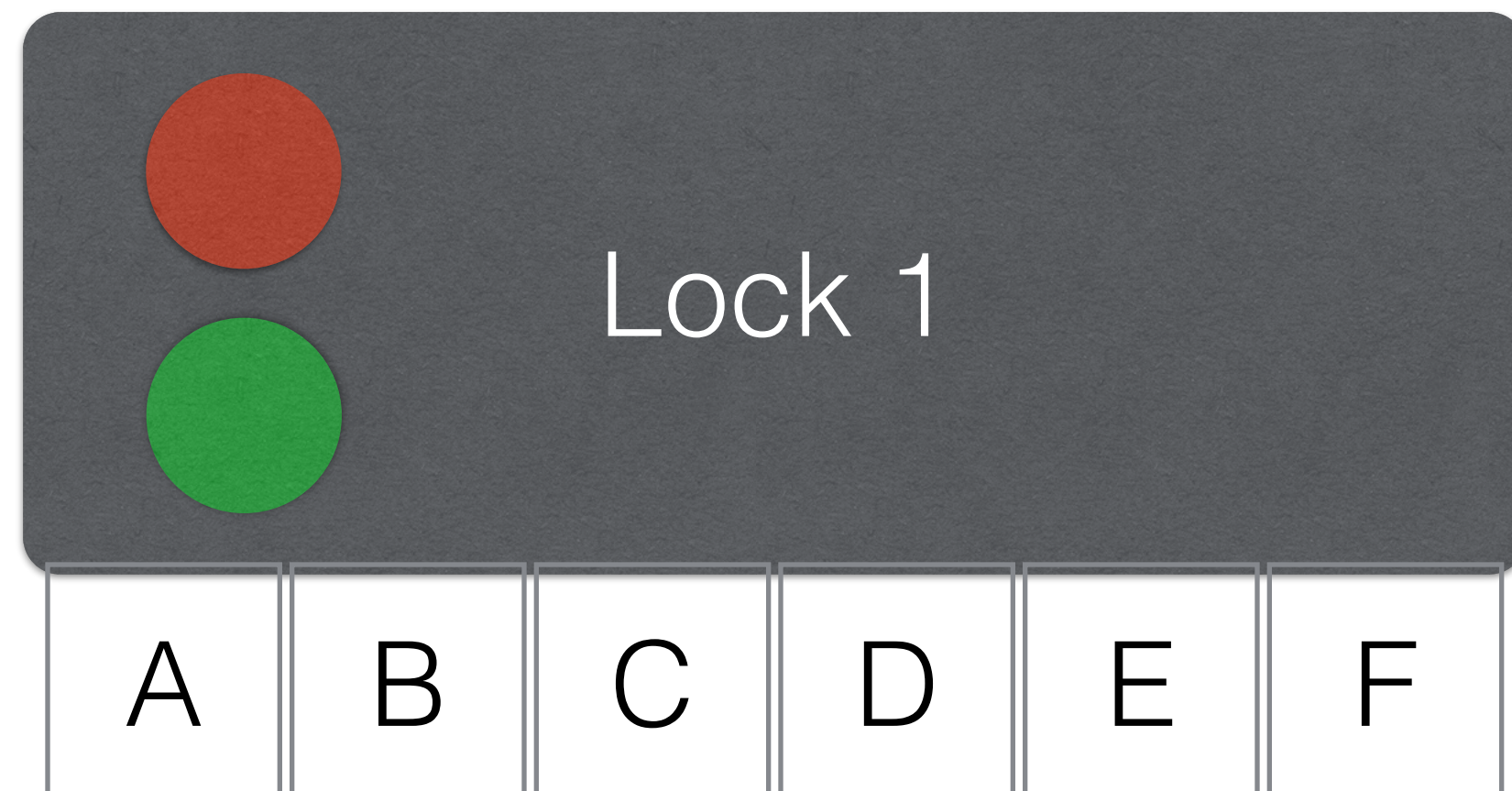
# Locks

- Coarse-grained locking

| A | B | C | D | E | F |

# Locks

- Coarse-grained locking

# Locks

- Coarse-grained locking



Lock 1

| A | B | C | D | E | F |

**TASK 1**

```
Get_Lock(1)
Write(A, x)
Release_Lock(1)
```

# Locks

- Coarse-grained locking



Lock 1

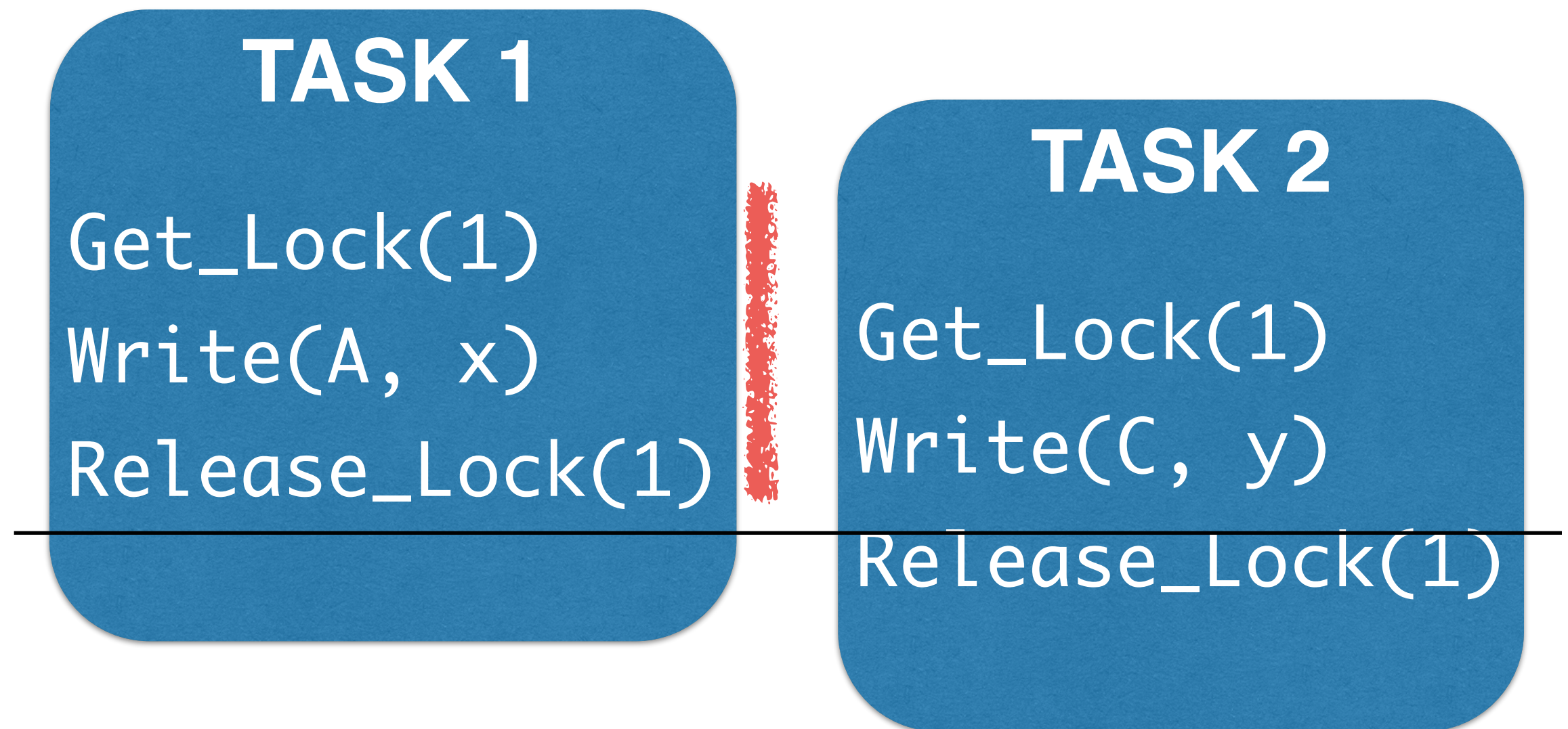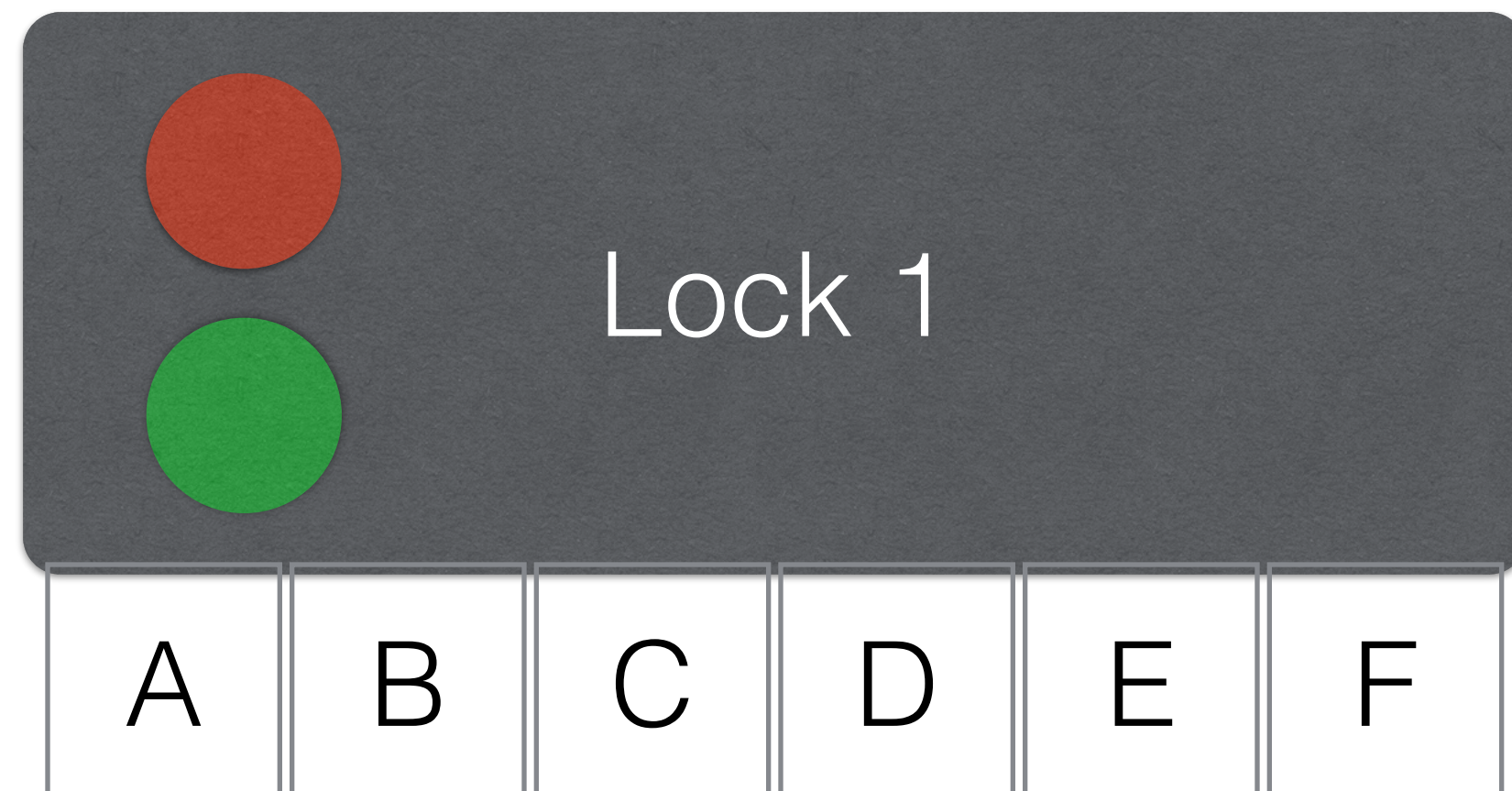| A | B | C | D | E | F |

**TASK 1**

```
Get_Lock(1)
Write(A, x)
Release_Lock(1)
```
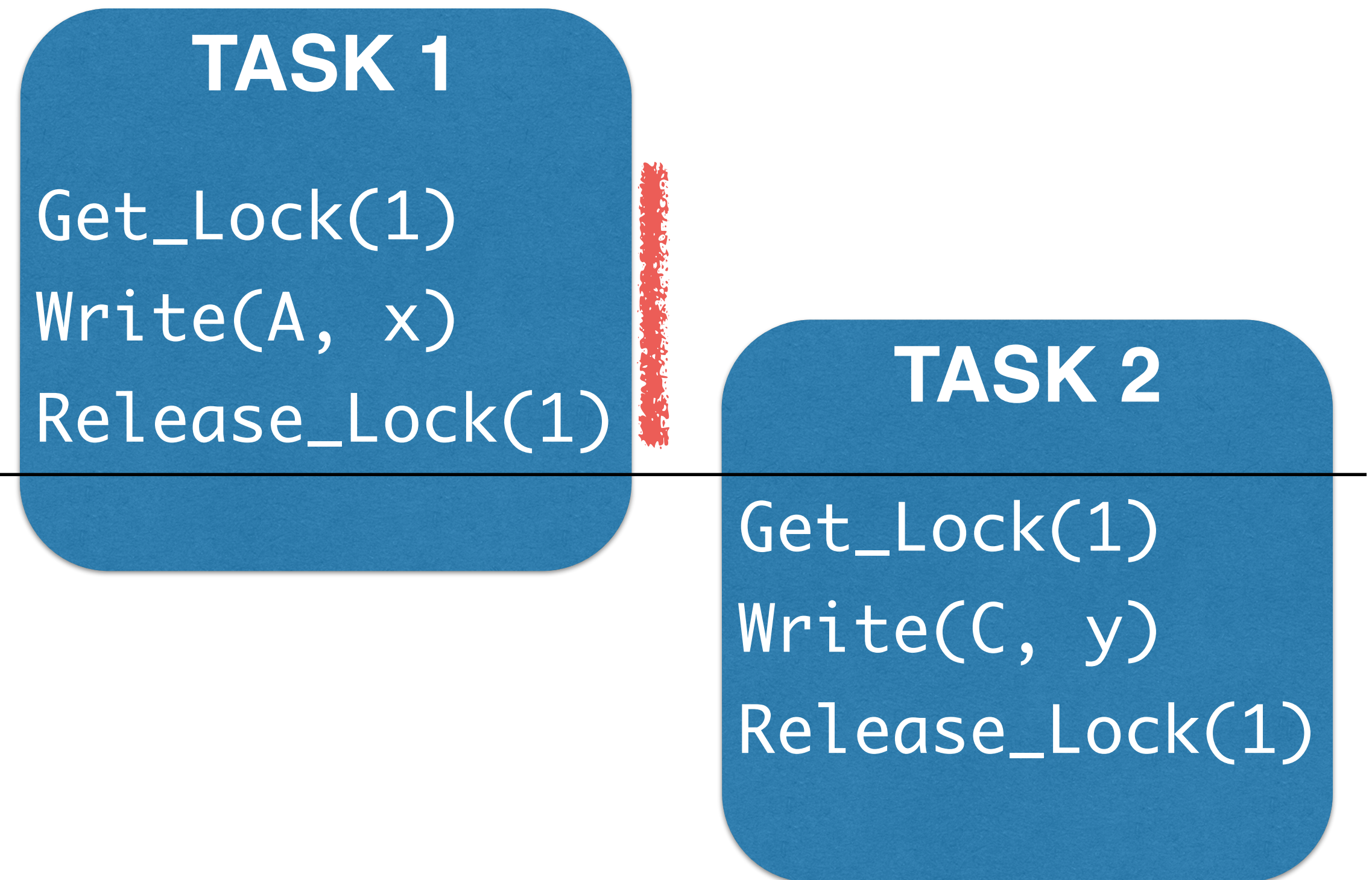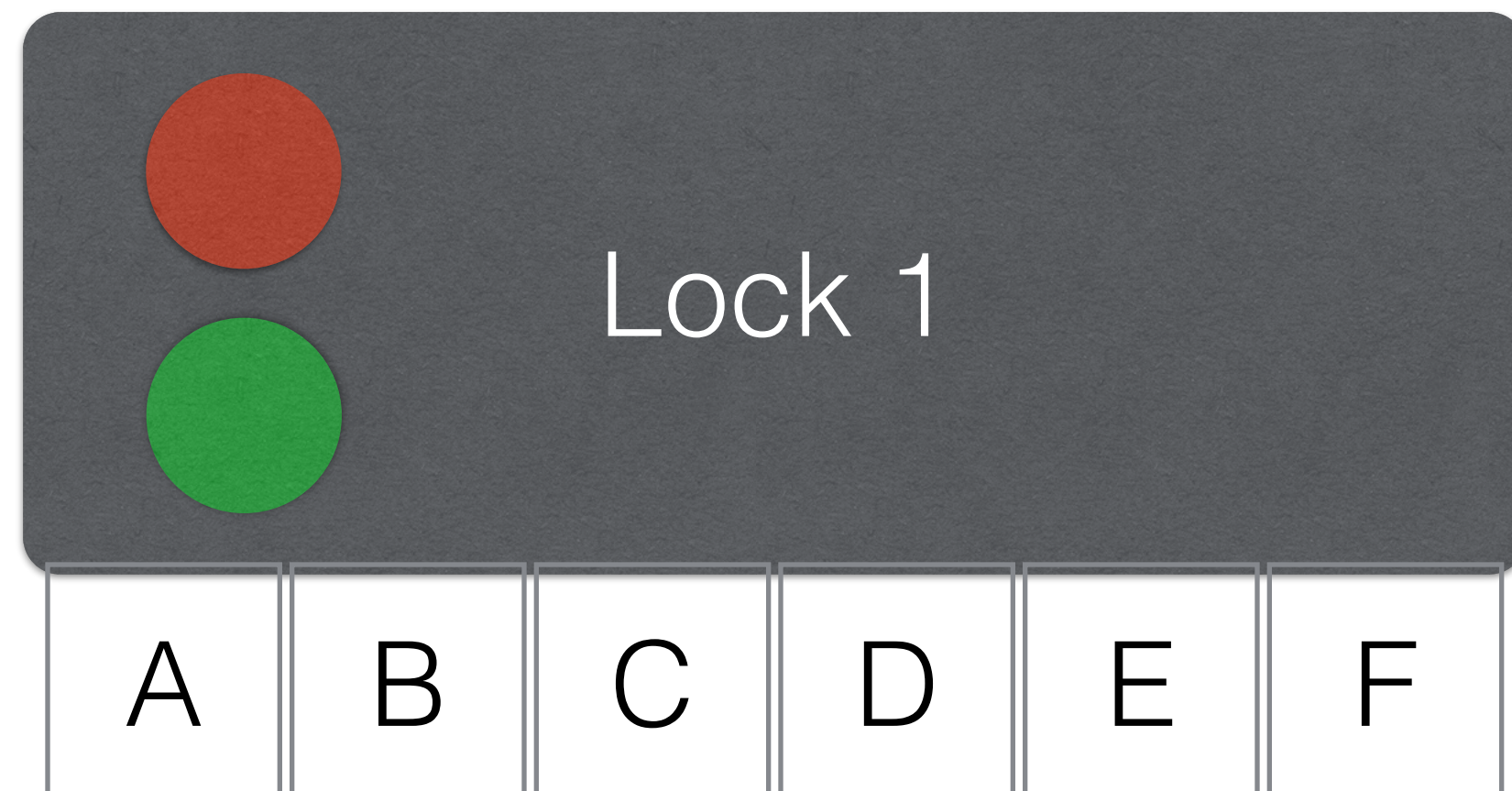
**TASK 2**

```
Get_Lock(1)
Write(C, y)
Release_Lock(1)
```

# Locks

- Coarse-grained locking

**Lock 1**

| A | B | C | D | E | F |

**TASK 1**

```
Get_Lock(1)
Write(A, x)
Release_Lock(1)
```

**TASK 2**

```
Get_Lock(1)
Write(C, y)
Release_Lock(1)
```
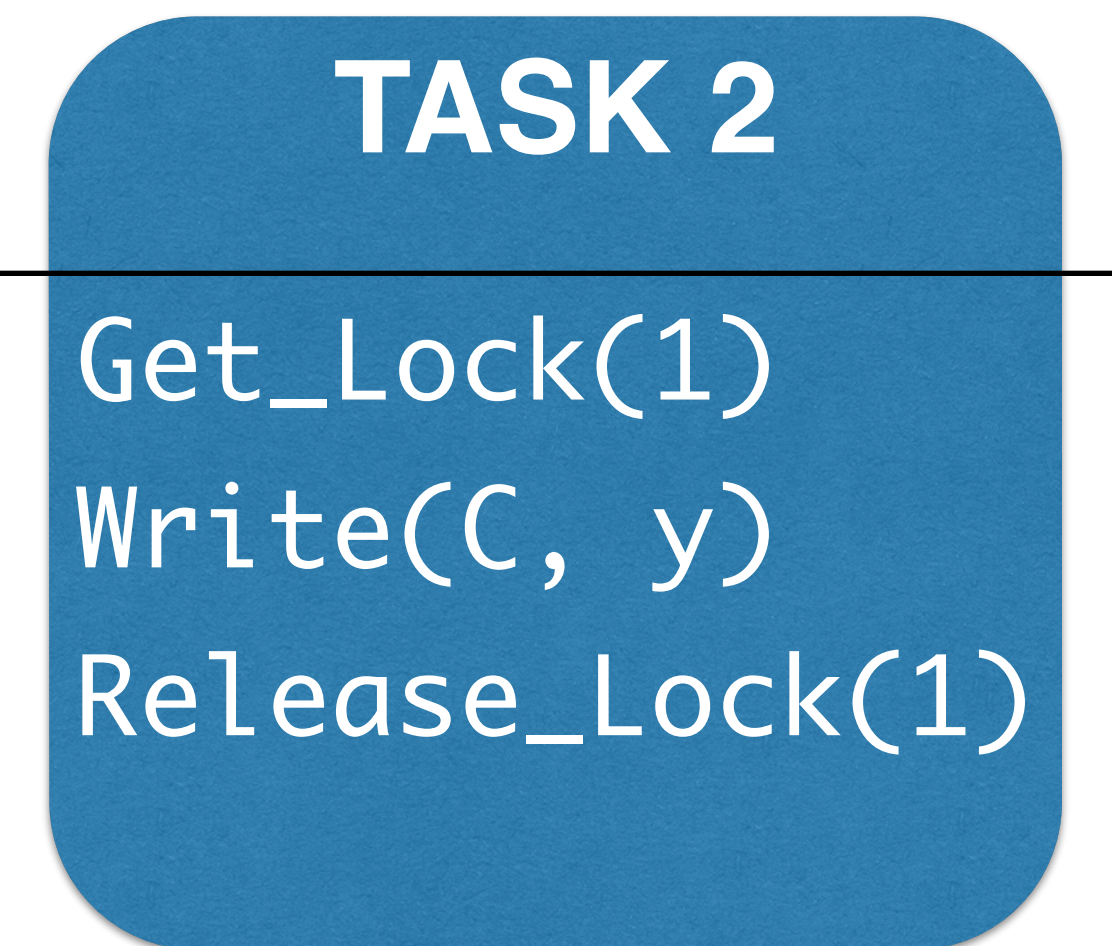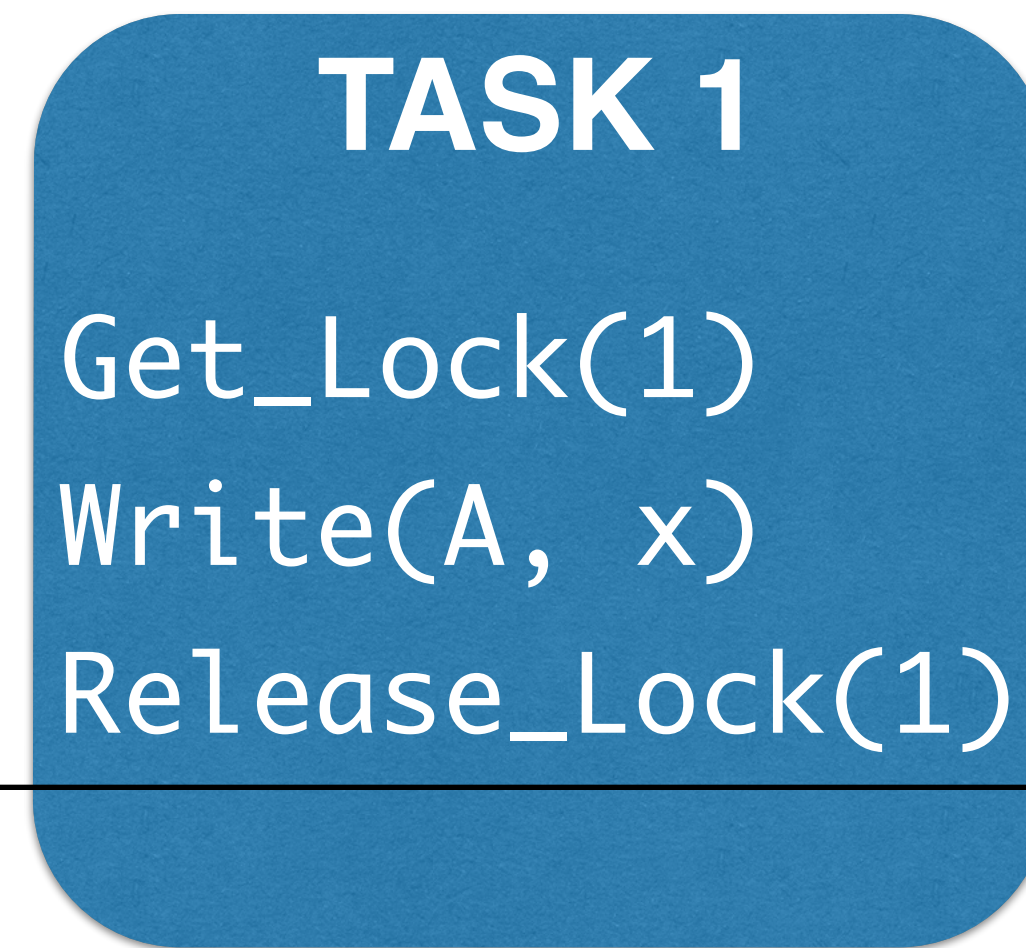
# Locks

- Coarse-grained locking



**Lock 1**

| A | B | C | D | E | F |

**TASK 1**

```
Get_Lock(1)
Write(A, x)
Release_Lock(1)
```

**TASK 2**

```
Get_Lock(1)
Write(C, y)
Release_Lock(1)
```

# Locks

- Coarse-grained locking

Lock 1

| A | B | C | D | E | F |

Critical sections can not progress in parallel!

**TASK 1**

```
Get_Lock(1)
Write(A, x)
Release_Lock(1)
```

**TASK 2**

```
Get_Lock(1)
Write(C, y)
Release_Lock(1)
```
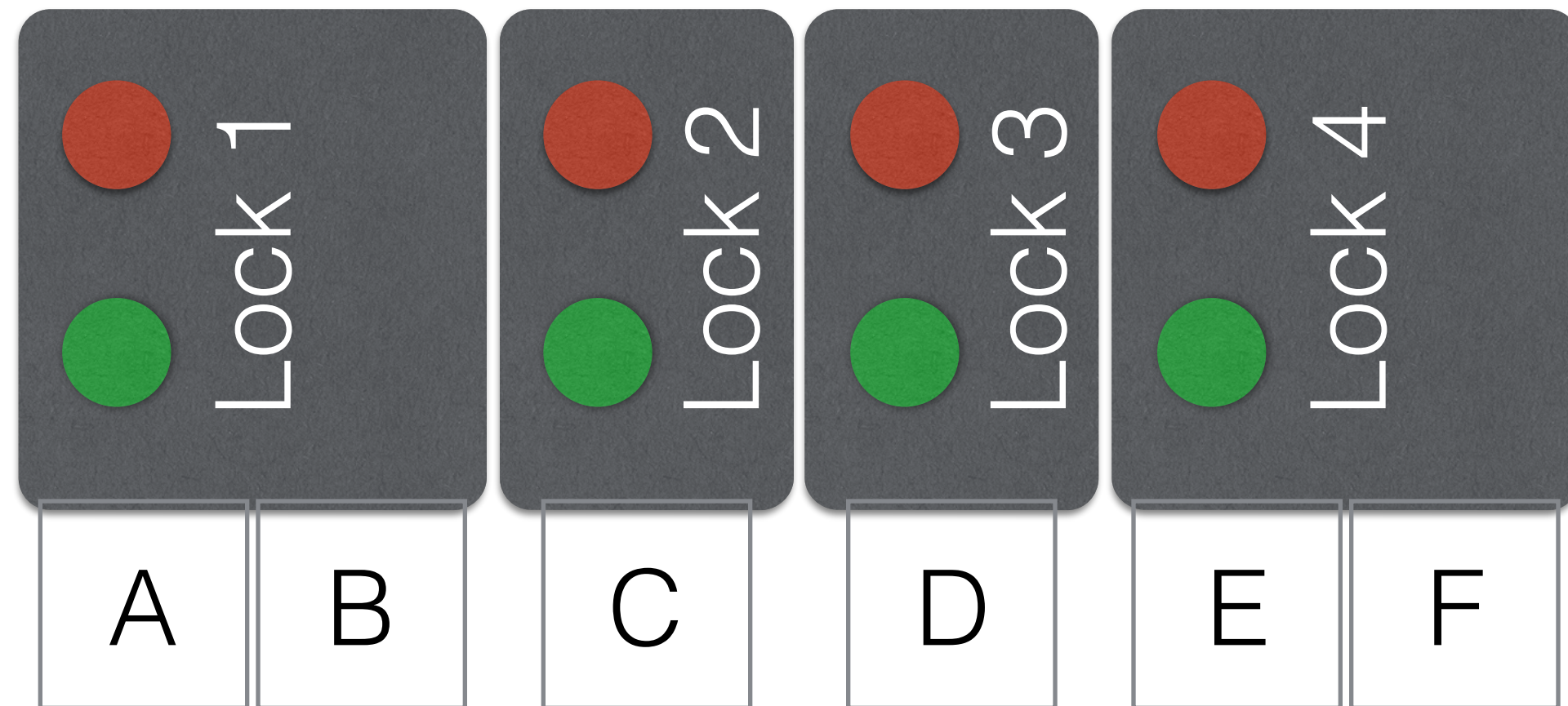
# Locks

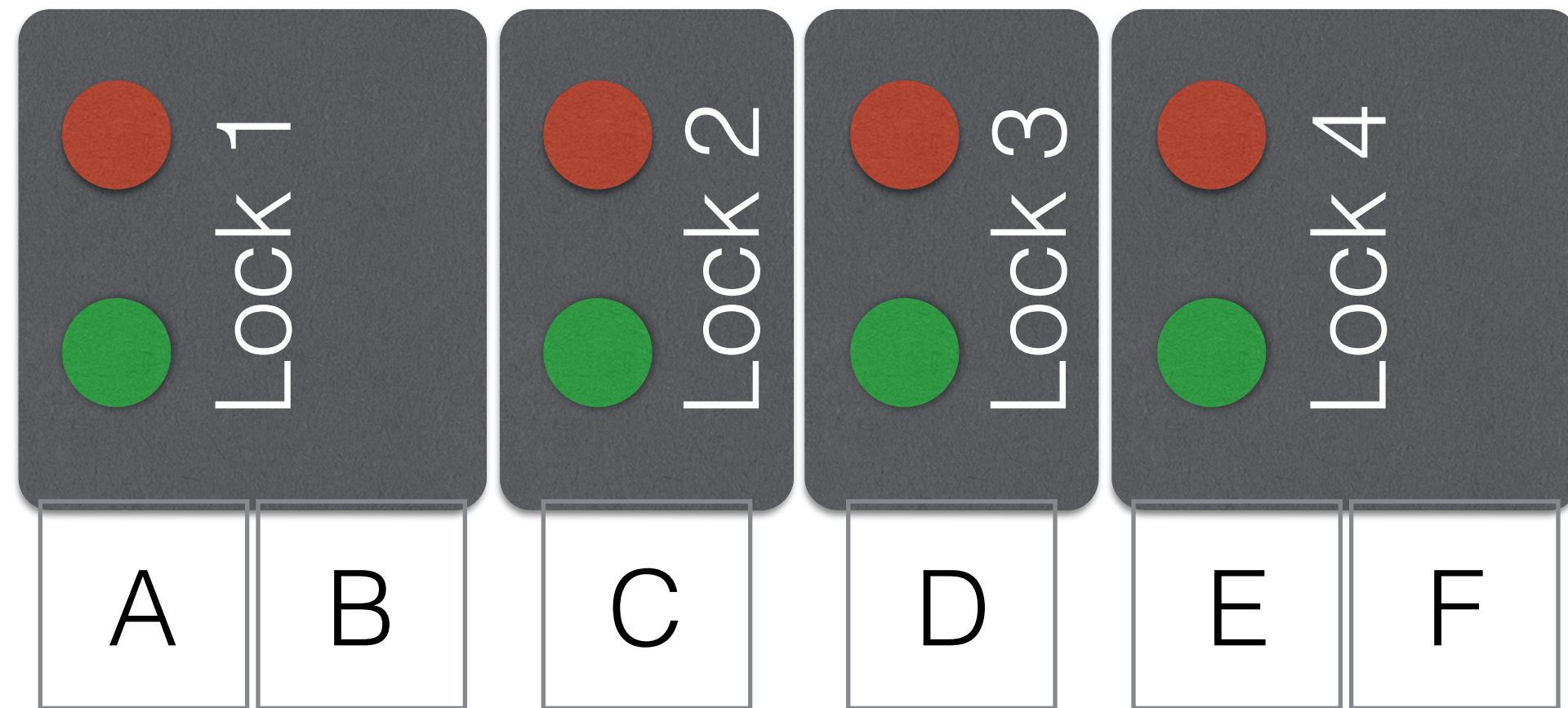- Fine-grained locking

# Locks

- Fine-grained locking

# Locks

- Fine-grained locking

# Locks

- Fine-grained locking



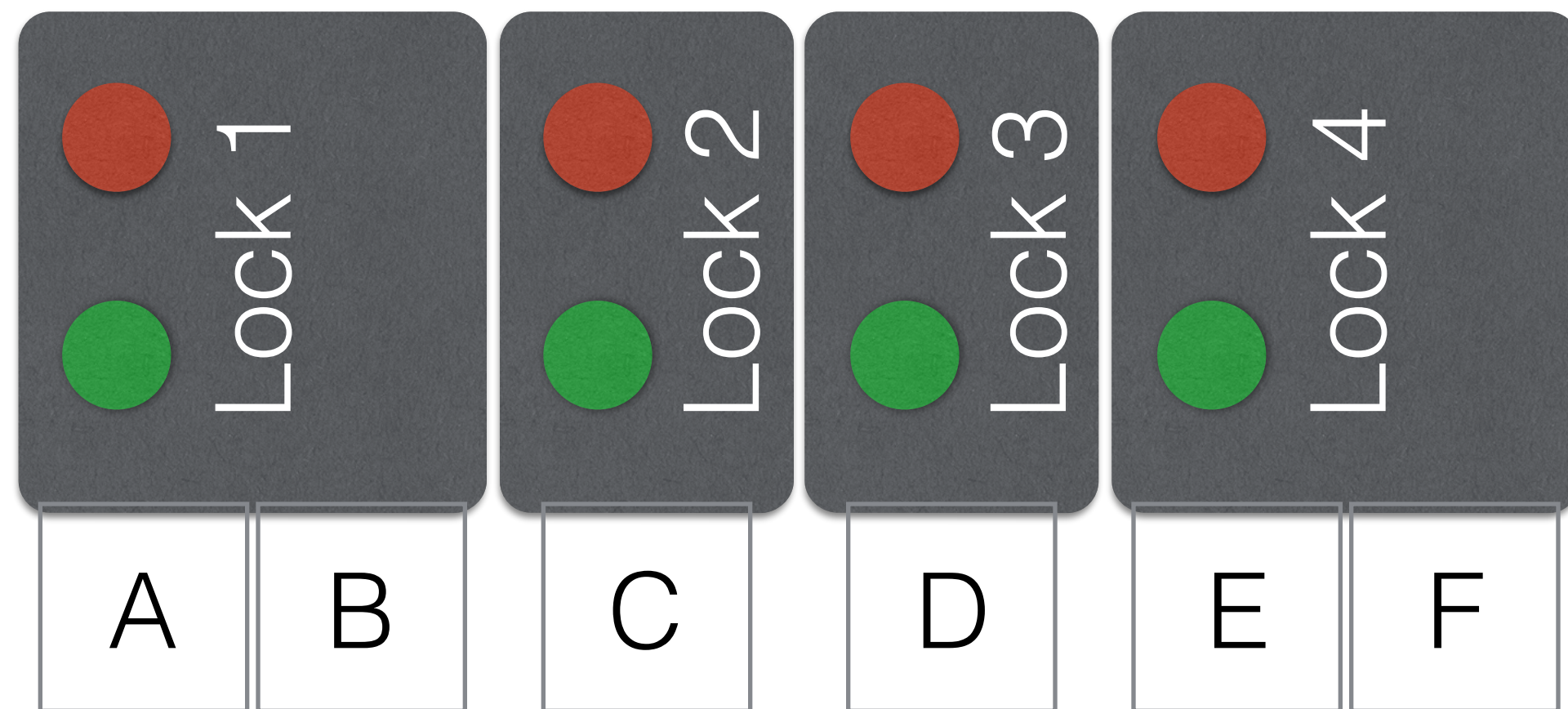**TASK 1**

```
Get_Lock(1)
Get_Lock(3)
Read(C, x)
Write(A, x)
Release_Lock(3)
Release_Lock(1)
```

# Locks

- Fine-grained locking



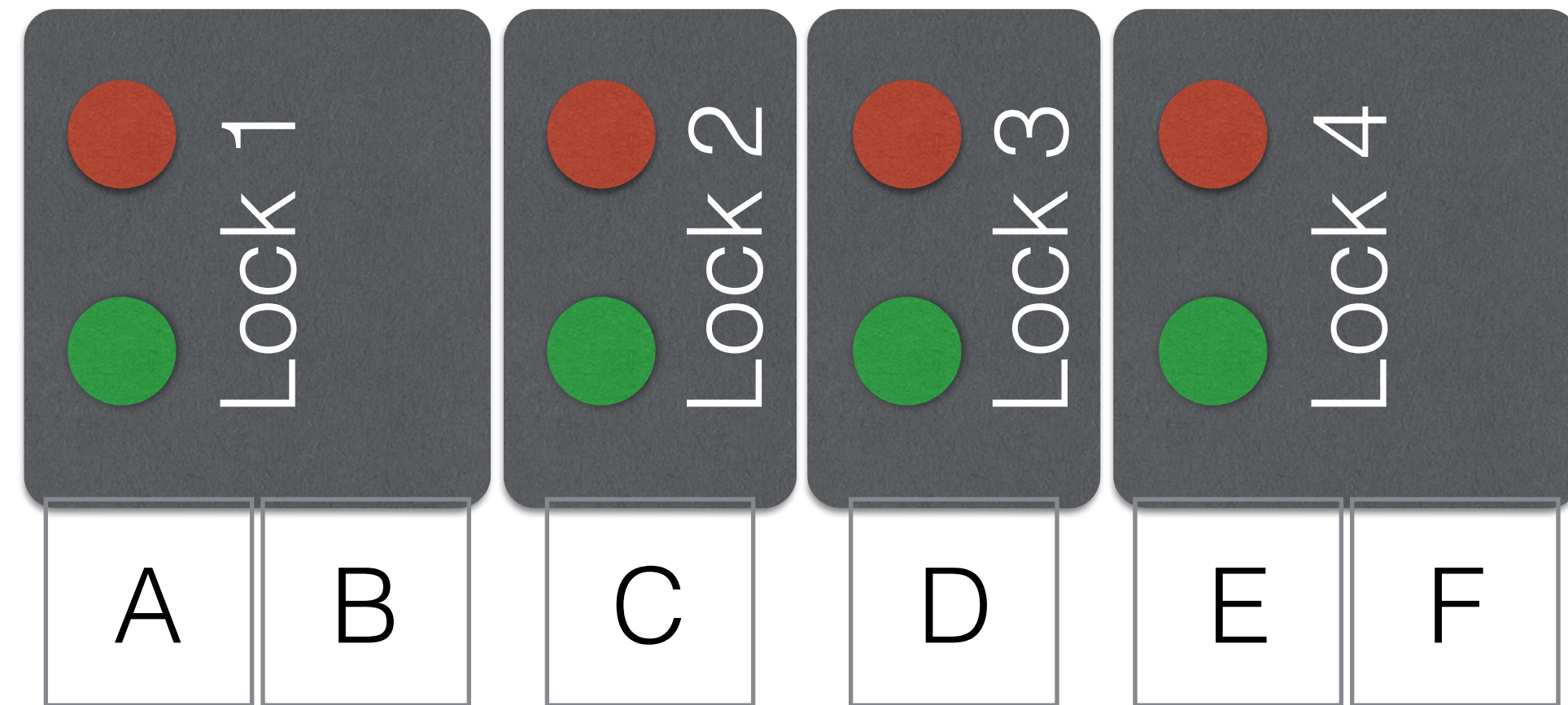| Lock 1 | Lock 2 | Lock 3 | Lock 4 |
|--------|--------|--------|--------|
| A B | C | D | E F |

**TASK 1**

```
Get_Lock(1)
Get_Lock(3)
Read(C, x)
Write(A, x)
Release_Lock(3)
Release_Lock(1)
```

**TASK 2**

```
Get_Lock(3)
Get_Lock(1)
Read(B, x)
Write(C, y)
Release_Lock(1)
Release_Lock(3)
```
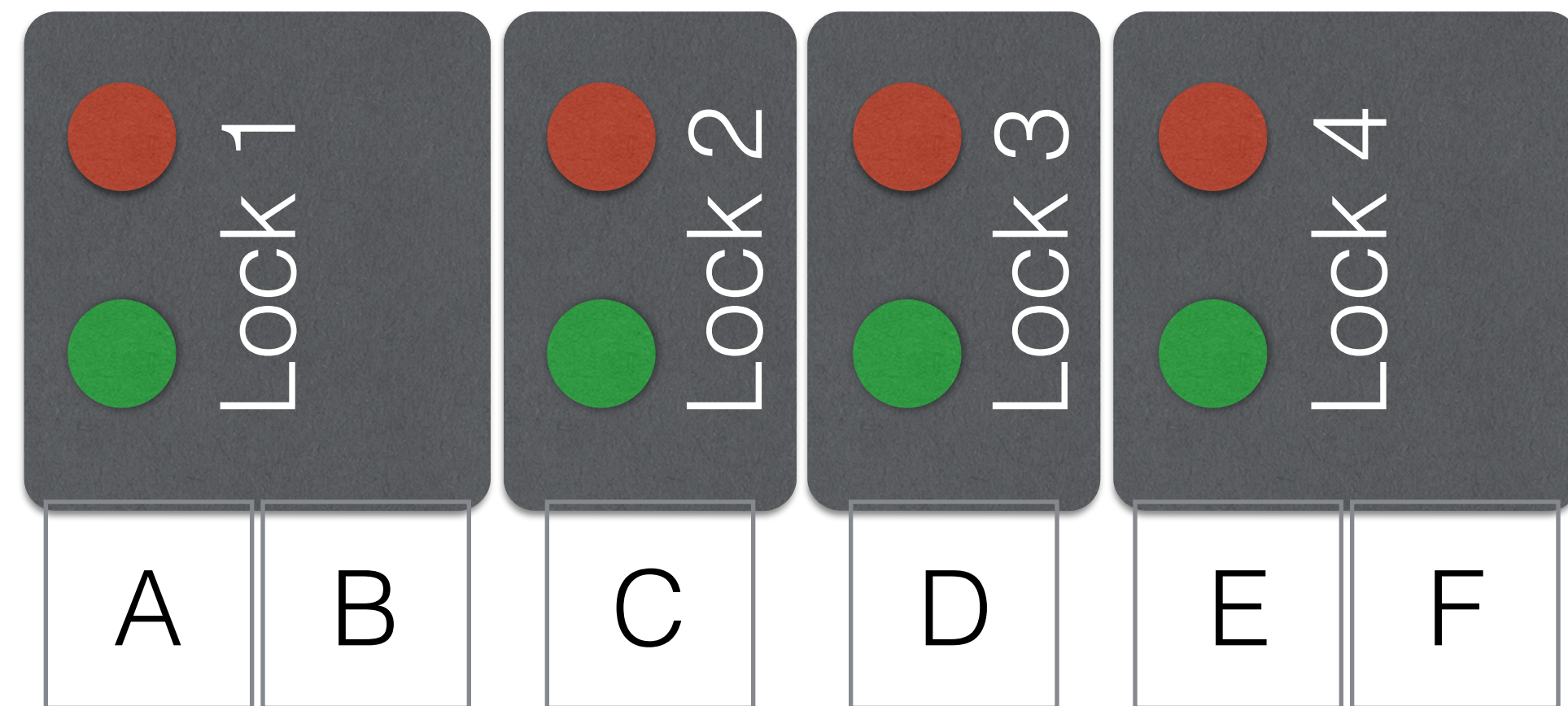
# Locks

- Fine-grained locking

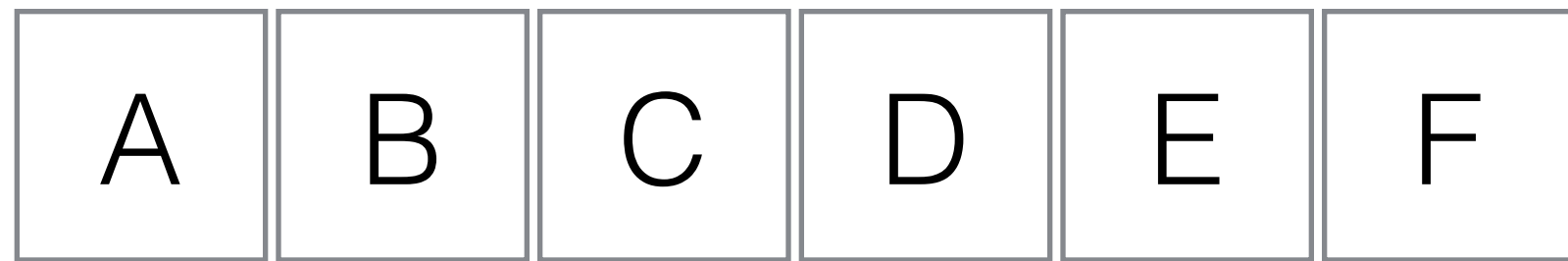| | | | | |
|---|---|---|---|---|
| 🔴 🟢 Lock 1 | 🔴 🟢 Lock 2 | 🔴 🟢 Lock 3 | 🔴 🟢 Lock 4 | |
| A | B | C | D | E | F |

**TASK 1**

```
Get_Lock(1)
Get_Lock(3)
Read(C, x)
Write(A, x)
Release_Lock(3)
Release_Lock(1)
```

**TASK 2**

```
Get_Lock(3)
Get_Lock(1)
Read(B, x)
Write(C, y)
Release_Lock(1)
Release_Lock(3)
```

# Locks

- Fine-grained locking



**TASK 1**

```
Get_Lock(1)
Get_Lock(3)
Read(C, x)
Write(A, x)
Release_Lock(3)
Release_Lock(1)
```

**TASK 2**

```
Get_Lock(3)
Get_Lock(1)
Read(B, x)
Write(C, y)
Release_Lock(1)
Release_Lock(3)
```
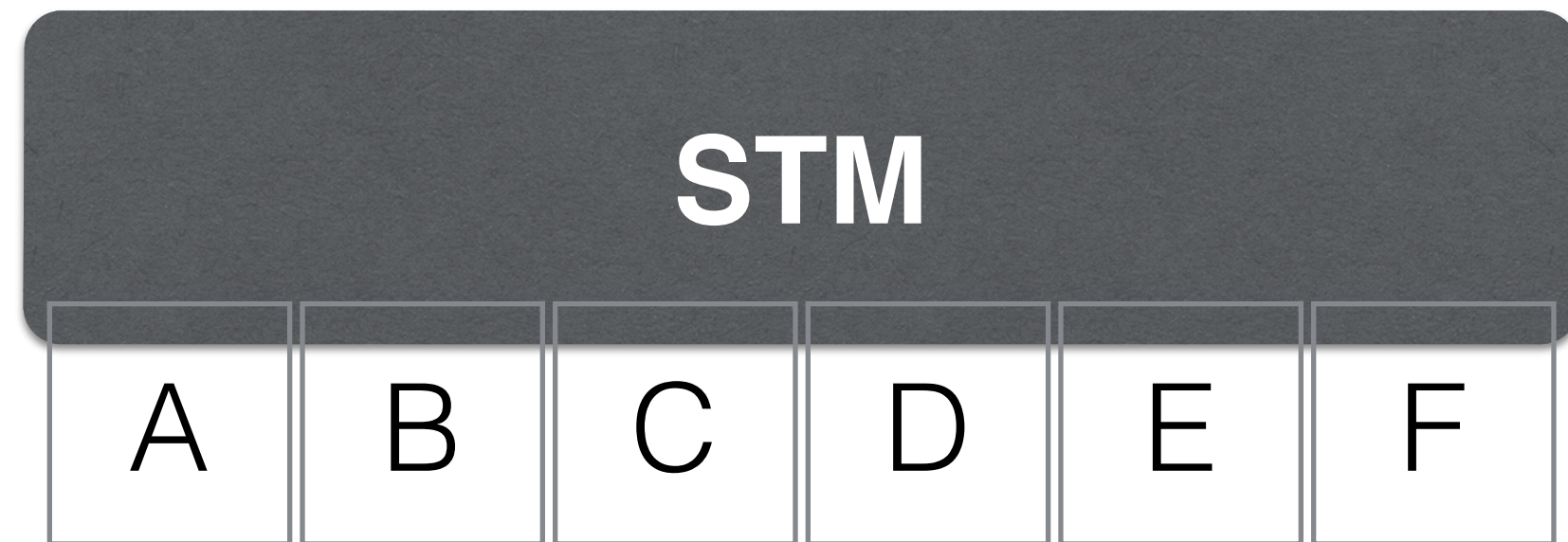
Increases system complexity with a negative impact on composability and maintainability!
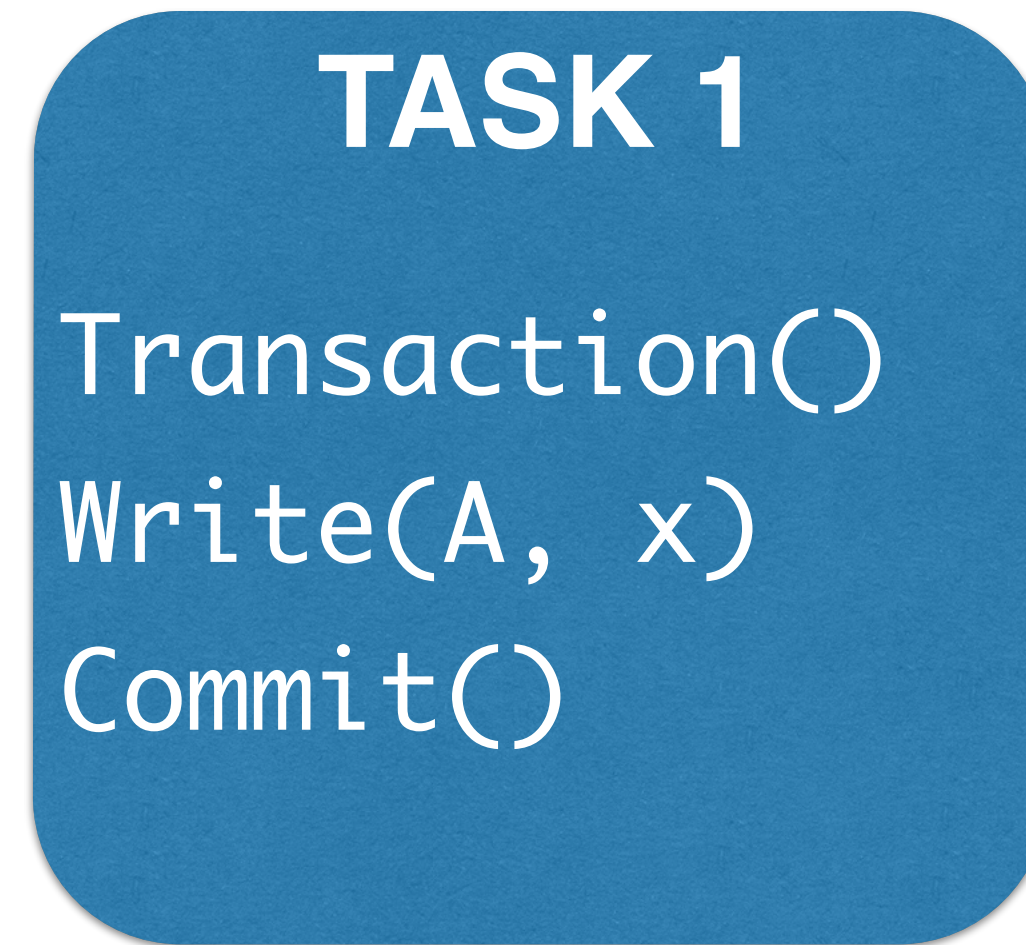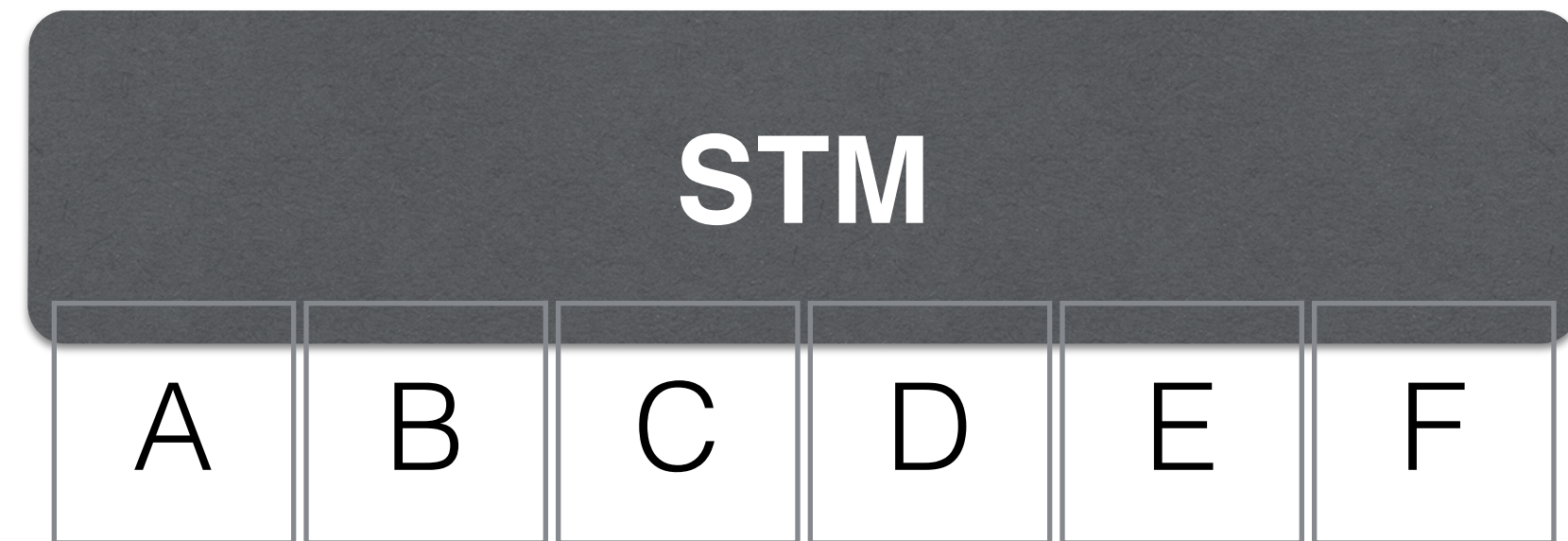
# Software Transactional Memory

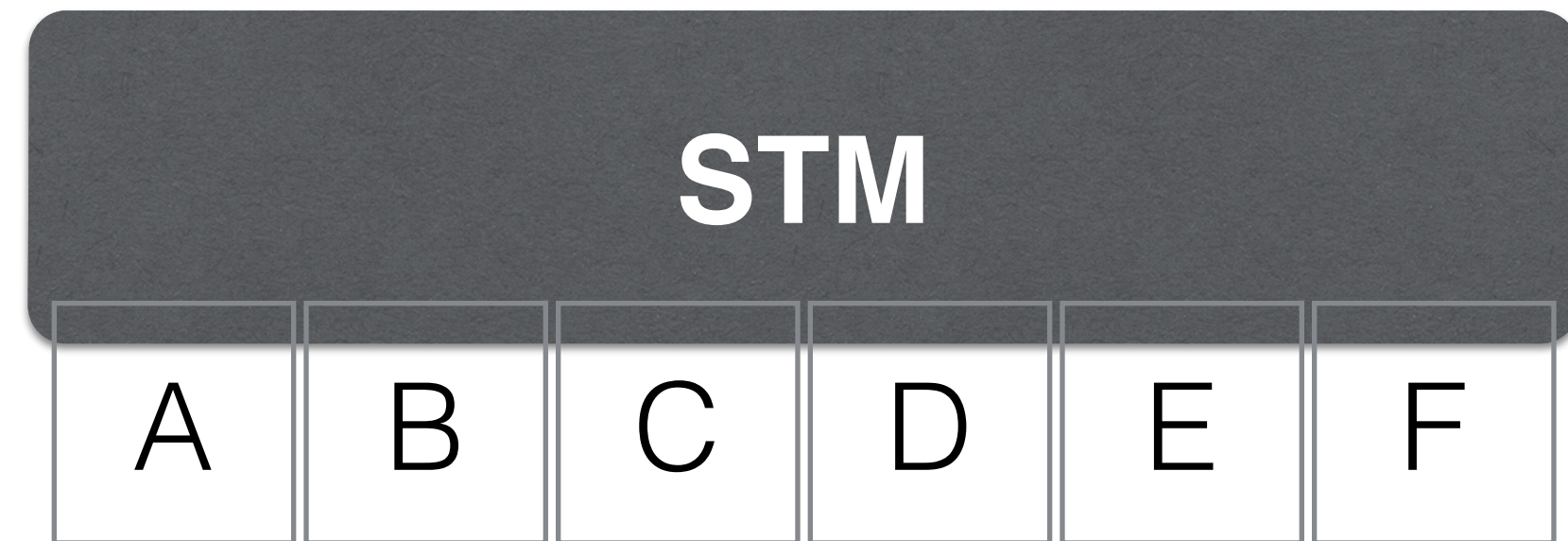| A | B | C | D | E | F |

# Software Transactional Memory

| STM | | | | | |
|---|---|---|---|---|---|
| A | B | C | D | E | F |

# Software Transactional Memory

**STM**

| A | B | C | D | E | F |

**TASK 1**

```
Transaction()
Write(A, x)
Commit()
```

# Software Transactional Memory

**STM**

| A | B | C | D | E | F |
|---|---|---|---|---|---|

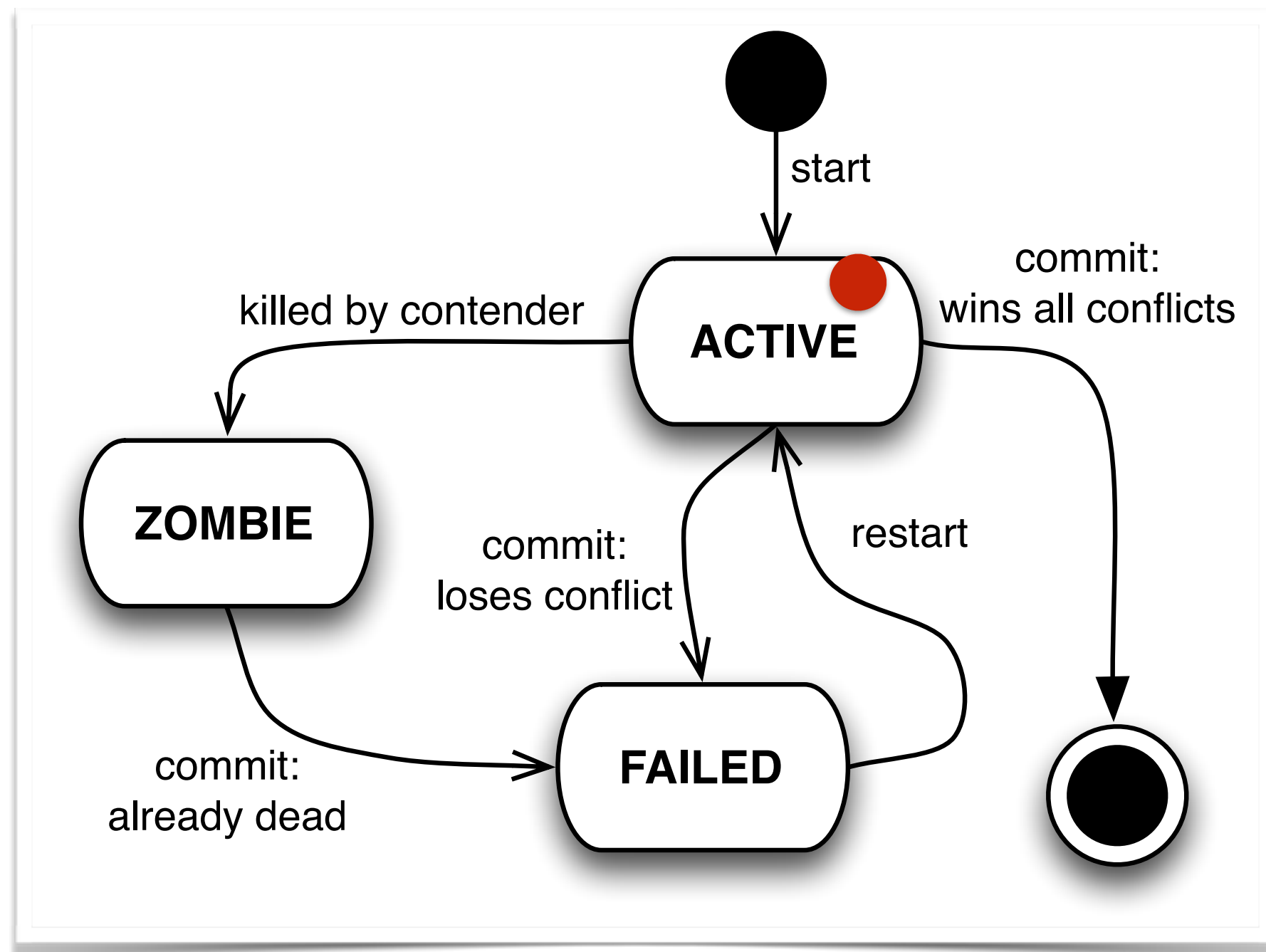**TASK 1**

```
Transaction()
Write(A, x)
Commit()
```

**TASK 2**

```
Transaction()
Write(C, y)
Commit()
```

# Software Transactional Memory



**STM**
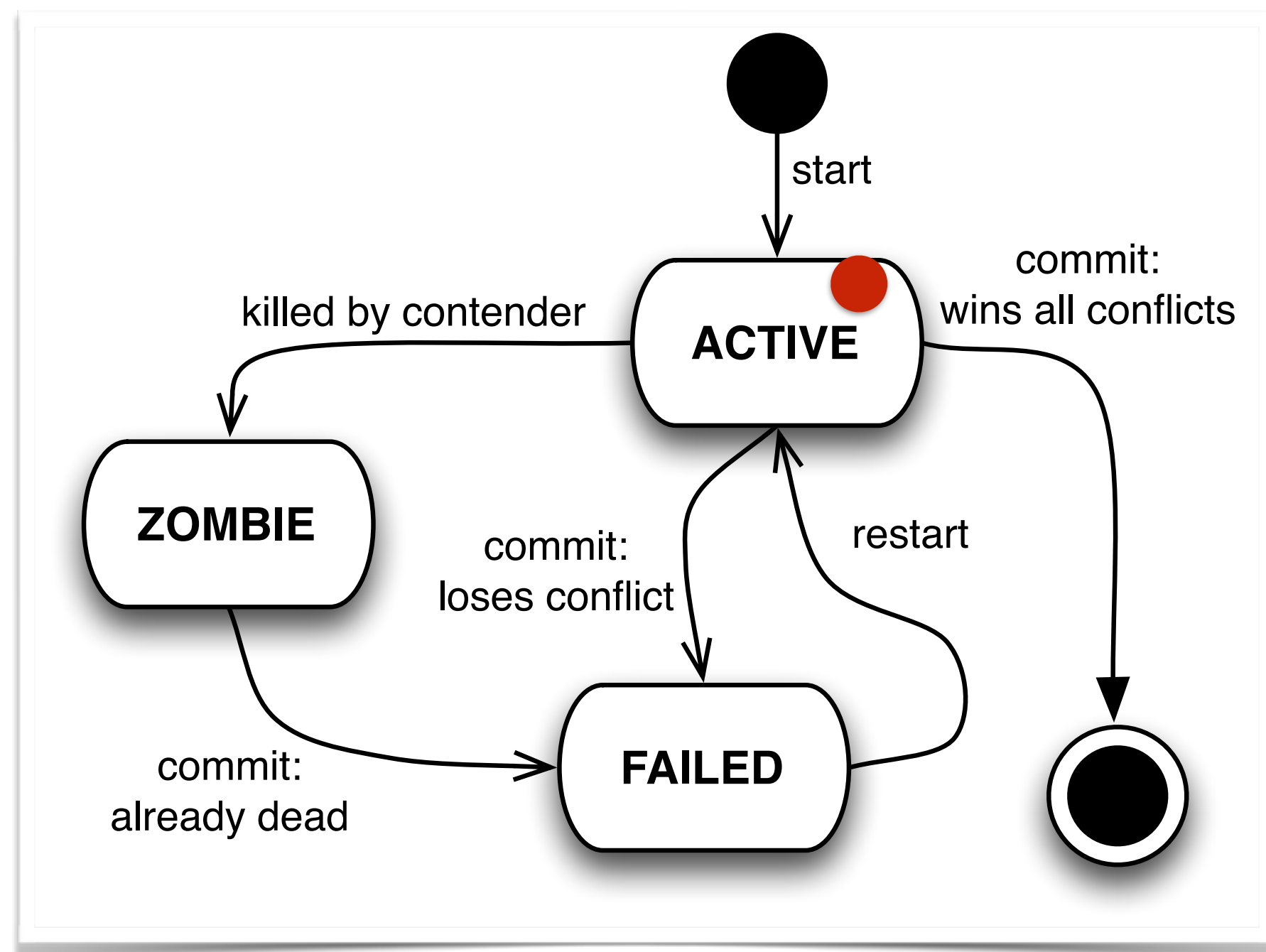
| A | B | C | D | E | F |

start

killed by contender

**ACTIVE**

commit:
wins all conflicts

**ZOMBIE**

commit:
loses conflict

restart

commit:
already dead

**FAILED**

**TASK 1**

```
Transaction()
Write(A, x)
Commit()
```

**TASK 2**

```
Transaction()
Write(C, y)
Commit()
```

# Software Transactional Memory

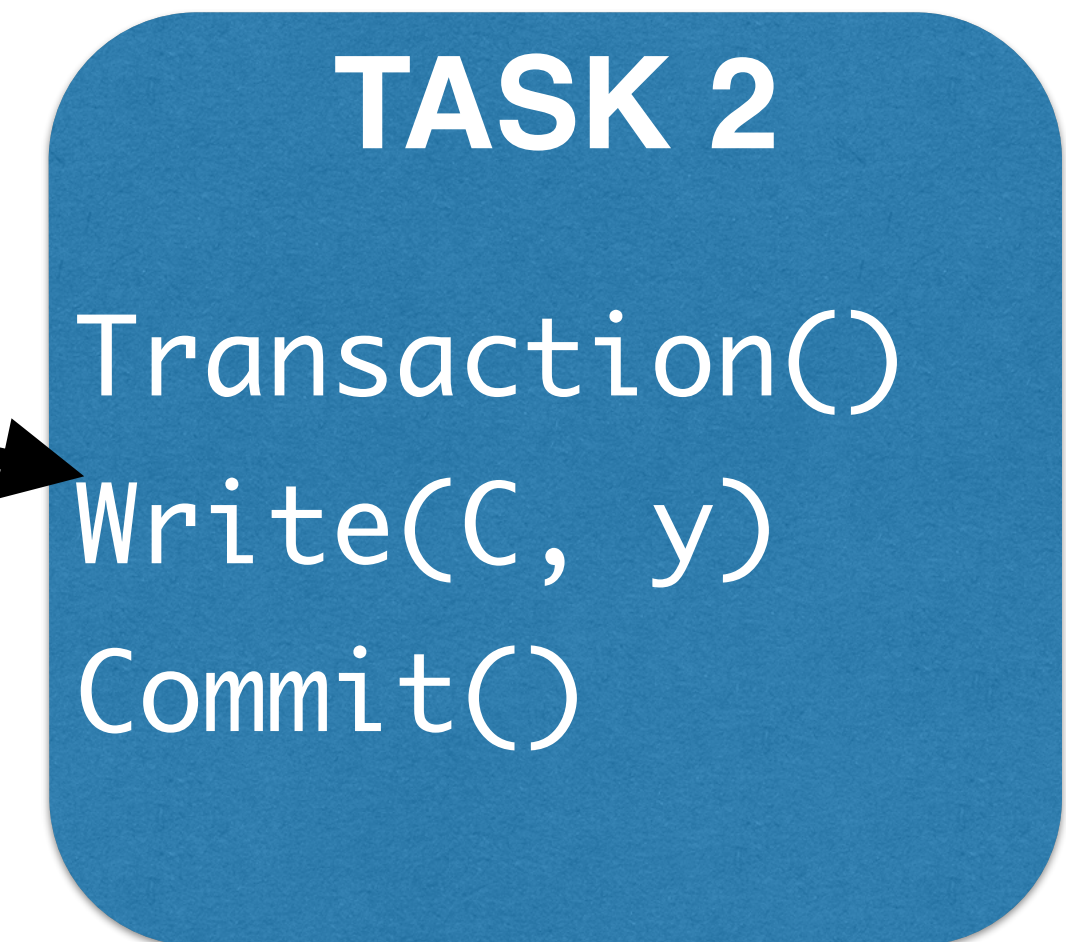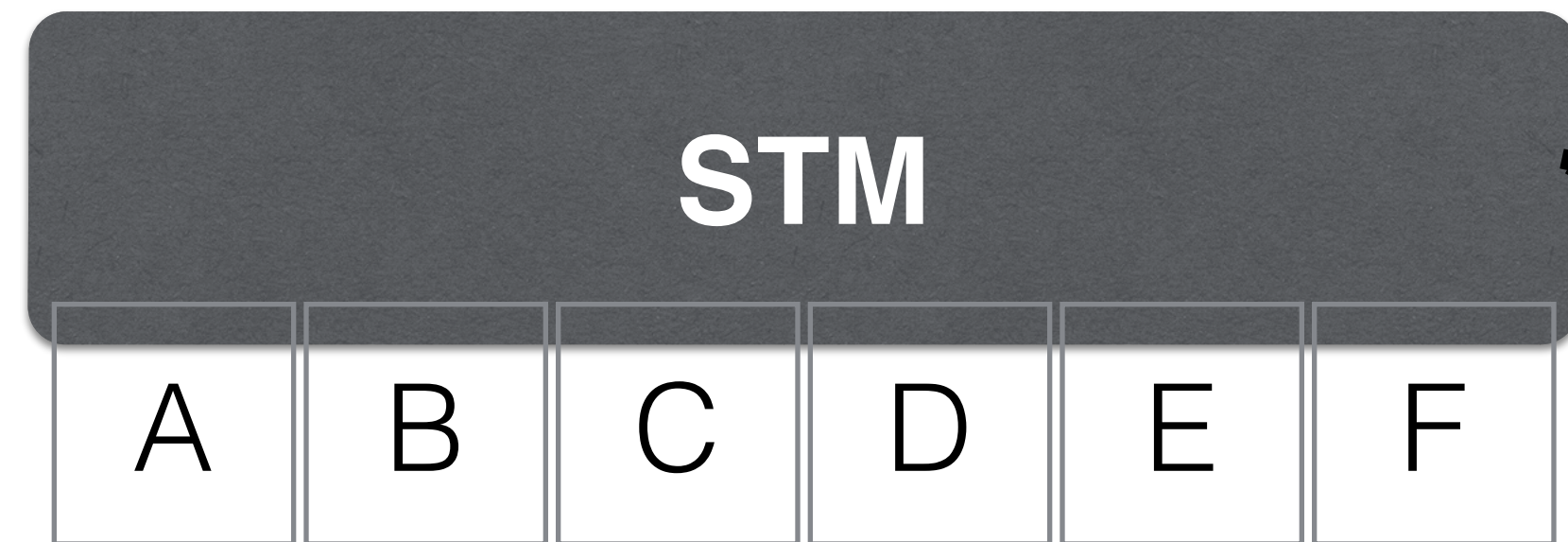# Software Transactional Memory

# Software Transactional Memory

# Software Transactional Memory



**STM**

| A | B | C | D | E | F |
|---|---|---|---|---|---|

**TASK 1**

```
Transaction()
Write(E, x)
Commit()
```

start

killed by contender

**ACTIVE**

commit:
wins all conflicts

**ZOMBIE**

commit:
loses conflict

restart

commit:
already dead

**FAILED**

# Software Transactional Memory



**STM**

| A | B | C | D | E | F |

start

ACTIVE

killed by contender

commit:
wins all conflicts

ZOMBIE

commit:
loses conflict

restart

commit:
already dead
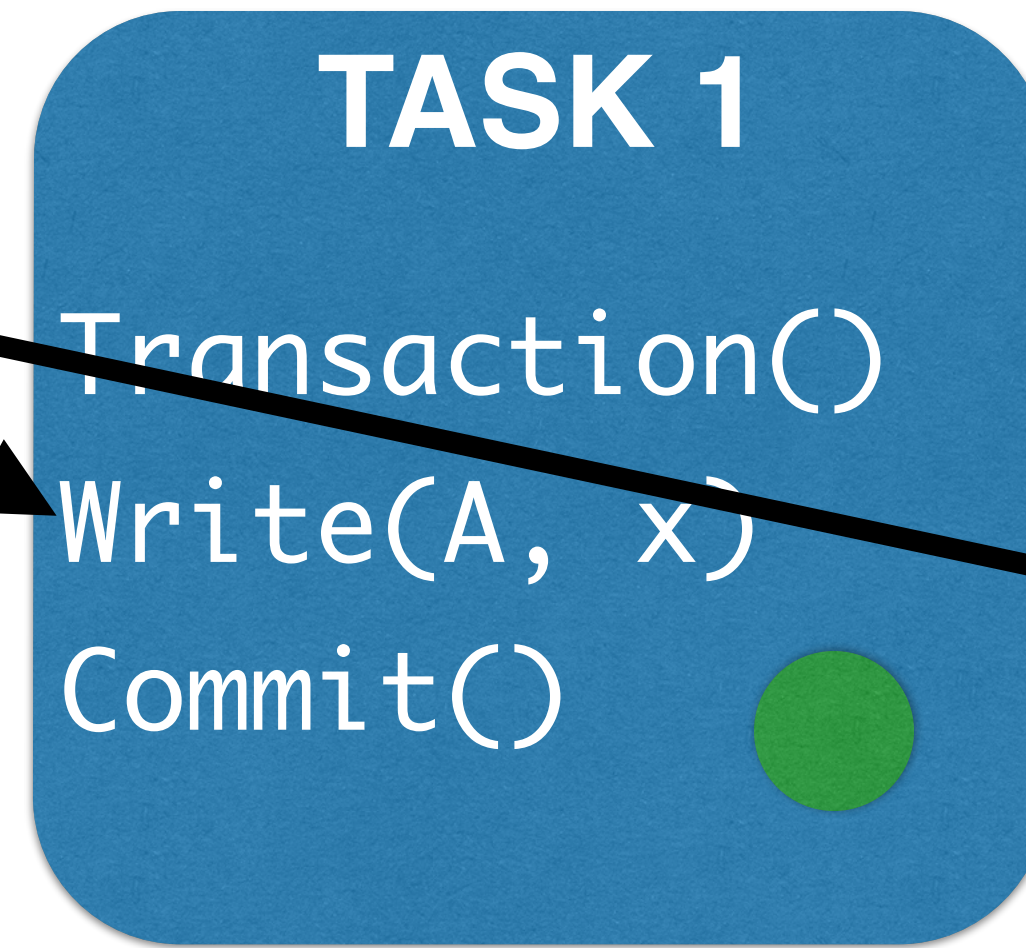
FAILED

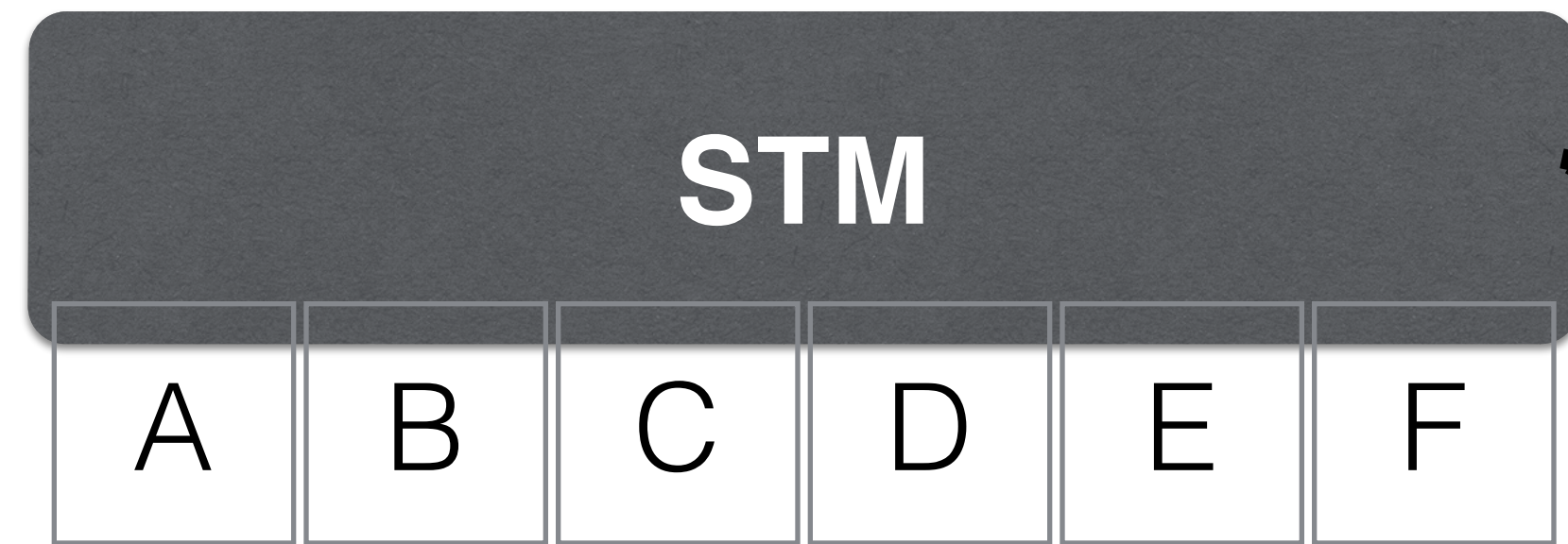**TASK 1**

```
Transaction()
Write(E, x)
Commit()
```

**TASK 2**

```
Transaction()
Write(E, y)
Commit()
```

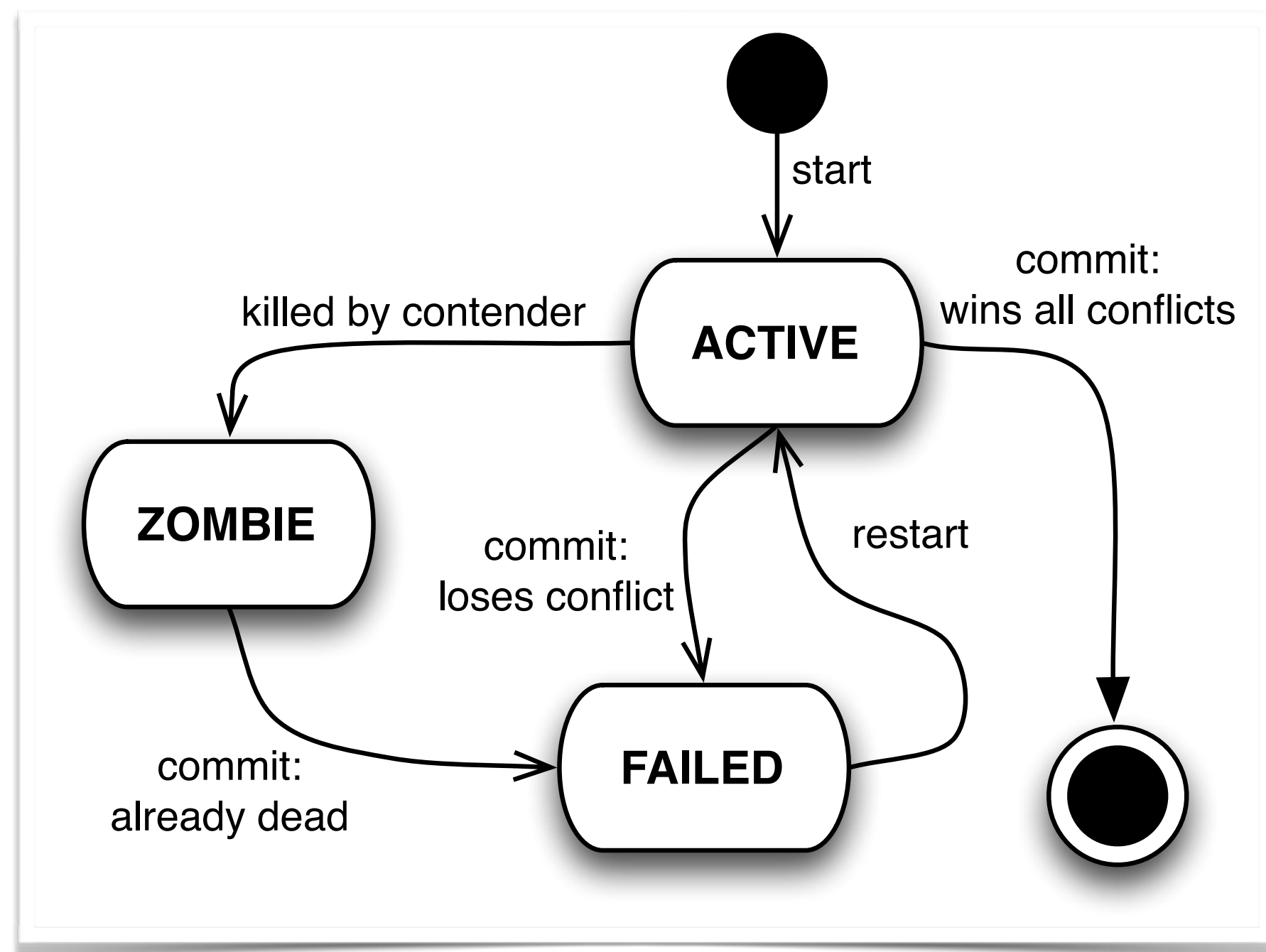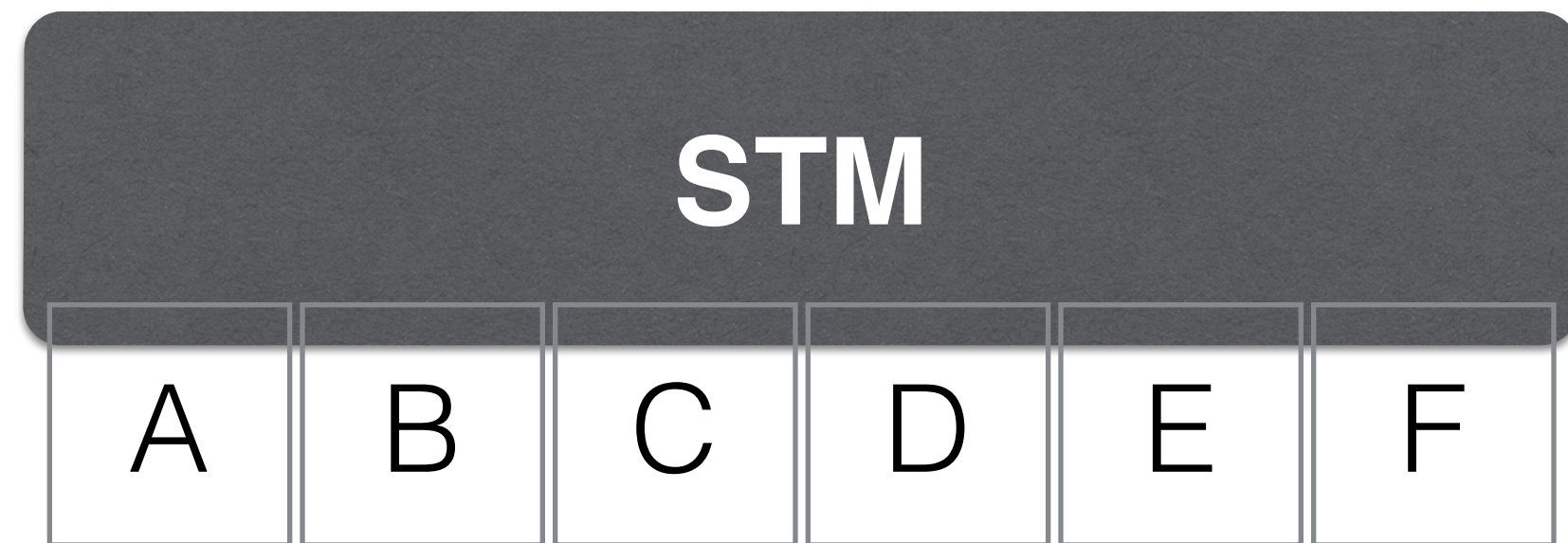# Software Transactional Memory

# Software Transactional Memory

**STM**
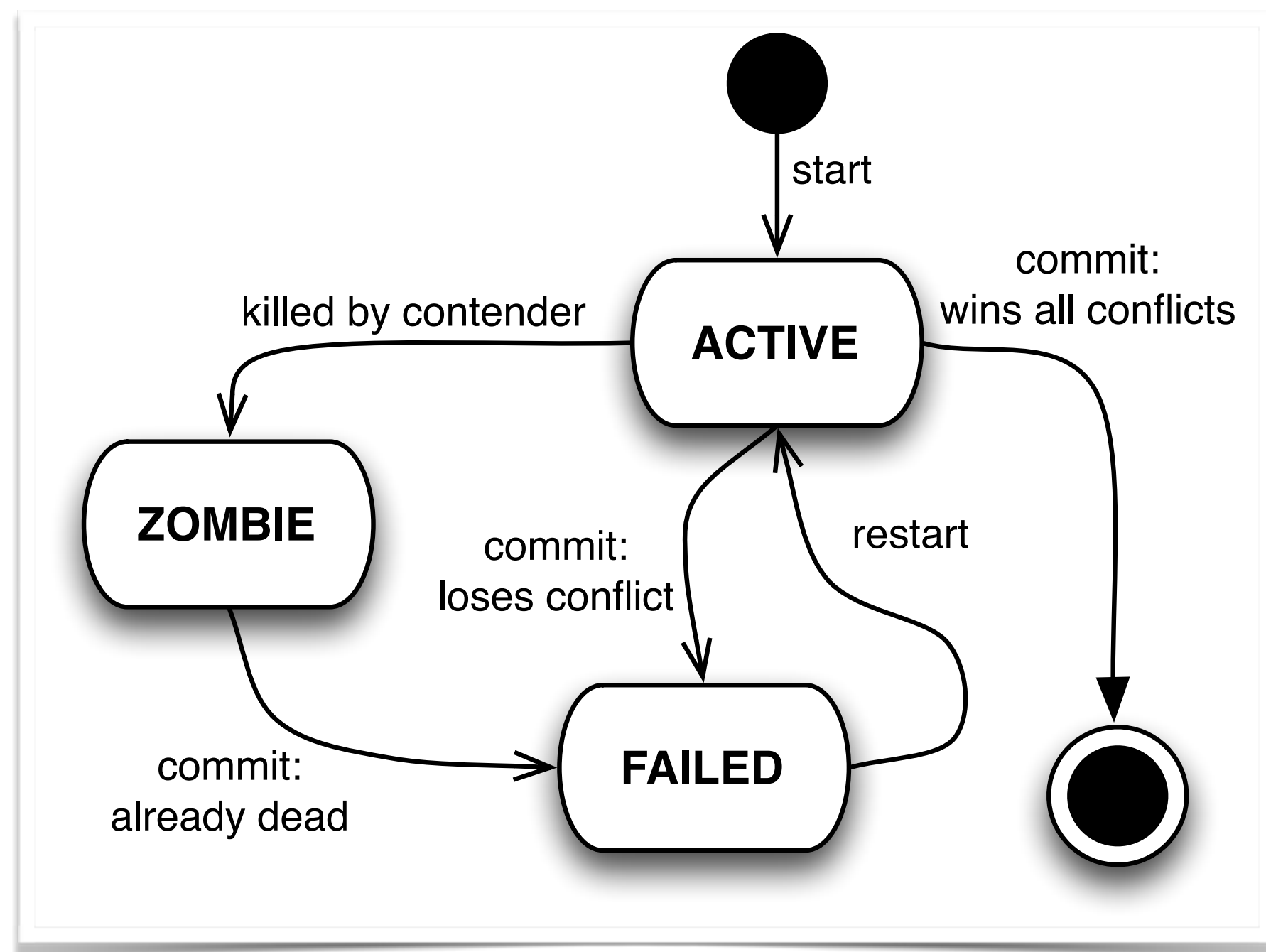
| A | B | C | D | E | F |
|---|---|---|---|---|---|

**TASK 1**

```
Transaction()
Write(E, x)
Commit()
```

**TASK 2**

```
Transaction()
Write(E, y)
Commit()
```

start

**ACTIVE**

commit: wins all conflicts

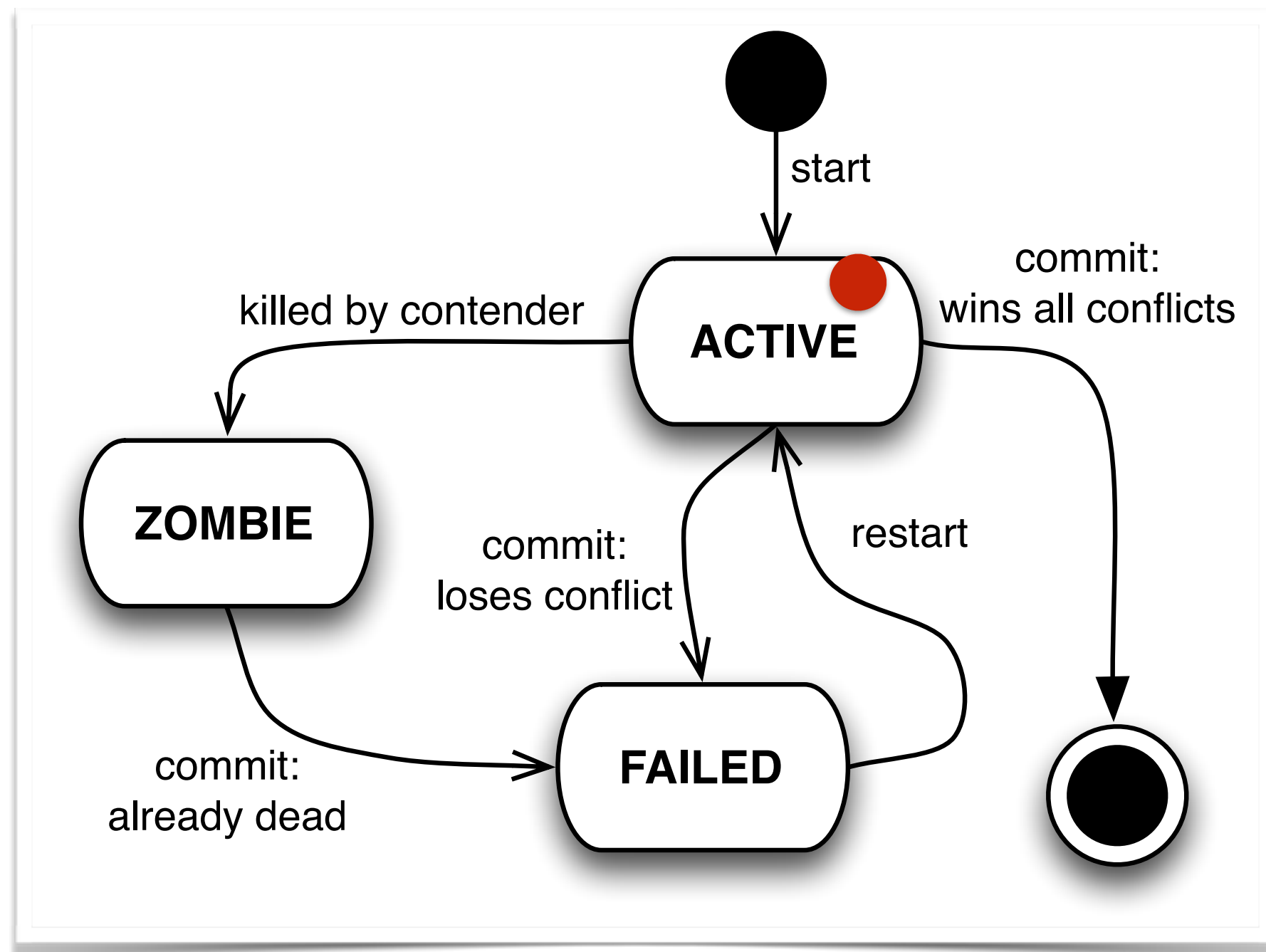killed by contender

**ZOMBIE**

commit: loses conflict

restart

commit: already dead

**FAILED**

# Software Transactional Memory

**STM**

| A | B | C | D | E | F |
|---|---|---|---|---|---|

**TASK 1**

```
Transaction()
Write(E, x)
Commit()
```

**TASK 2**

```
Transaction()
Write(E, y)
Commit()
```

**TASK 2**

```
Transaction()
Write(E, y)
Commit()
```

start

killed by contender

**ACTIVE**

commit:
wins all conflicts

**ZOMBIE**

commit:
loses conflict

restart

commit:
already dead

**FAILED**

# Software Transactional Memory

# Software Transactional Memory

# Managing contention

# Managing contention



**Polite**
Exponential back-off,
eventually commit.

# Managing contention



**Aggressive**
Kill the enemy!!!

# Managing contention



**Polite**
Exponential back-off, eventually commit.

**Aggressive**
Kill the enemy!!!

**Randomized**
Abort with p or Wait with (1-p).

# Managing contention

**Timestamp**
Older transaction survives.

**gressive**
e enemy!!!

# Managing contention



**Polite**
Exponential back-off, eventually commit.

**Aggressive**
Kill the enemy!!!

**Randomized**
Abort with p or
Wait with (1-p).

**Karma**
Accesses and aborts
accounts for karma.

# Managing contention



**Aggressive**
Kill the enemy!!!

...**mized**
...ith p or
...th (1-p).

**Eruption**
Priority rises if others are waiting.

# Managing contention

**DETERMINISTIC**

**NOT DETERMINISTIC**

**Timestamp**
Older transaction survives.

**Polite**
Exponential back-off, eventually commit.

**Aggressive**
Kill the enemy!!!

**Randomized**
Abort with p or Wait with (1-p).

**Karma**
Accesses and aborts accounts for karma.

**Eruption**
Priority rises if others are waiting.

# Model of computation and scheduling strategy

# Computation platform

- Multi-core

- Single memory bus

- Data shared in globally accessed memory, controlled by a STM system

# Application characteristics

- Application functionality divided into tasks.

- Each task is statically assigned to a core, before run-time.

- Each task releases a potentially infinite number of jobs.

  - Task: C (execution time), T (period), D (deadline)

  - Job: r (release time), d (absolute deadline)

T

D

C

r1    d1    r2    d2

# Application characteristics

- Application functionality divided into tasks.

- Each task is statically assigned to a core, before run-time.

- Each task releases a potentially infinite number of jobs.

  - Task: C (execution time), T (period), D (deadline)

  - Job: r (release time), d (absolute deadline)

# Serialisation of transactions in a RT environment

# FIFO serialisation of transactions

# FIFO serialisation of transactions

Problem solved!

- The order of serialisation of transactions in progress is determined once a transaction starts!

# FIFO serialisation of transactions

Problem solved!

- The order of serialisation of transactions in progress is determined once a transaction starts!

… or maybe not!

# FIFO serialisation of transactions

Problem solved!

- The order of serialisation of transactions in progress is determined once a transaction starts!

… or maybe not!

- What if jobs can be preempted while executing a transaction?

# FIFO serialisation of transactions

Problem solved!

- The order of serialisation of transactions in progress is determined once a transaction starts!

… or maybe not!

- What if jobs can be preempted while executing a transaction?

- What if multiple transactions can be simultaneously in progress on the same core?

# Preemptions and serialisation
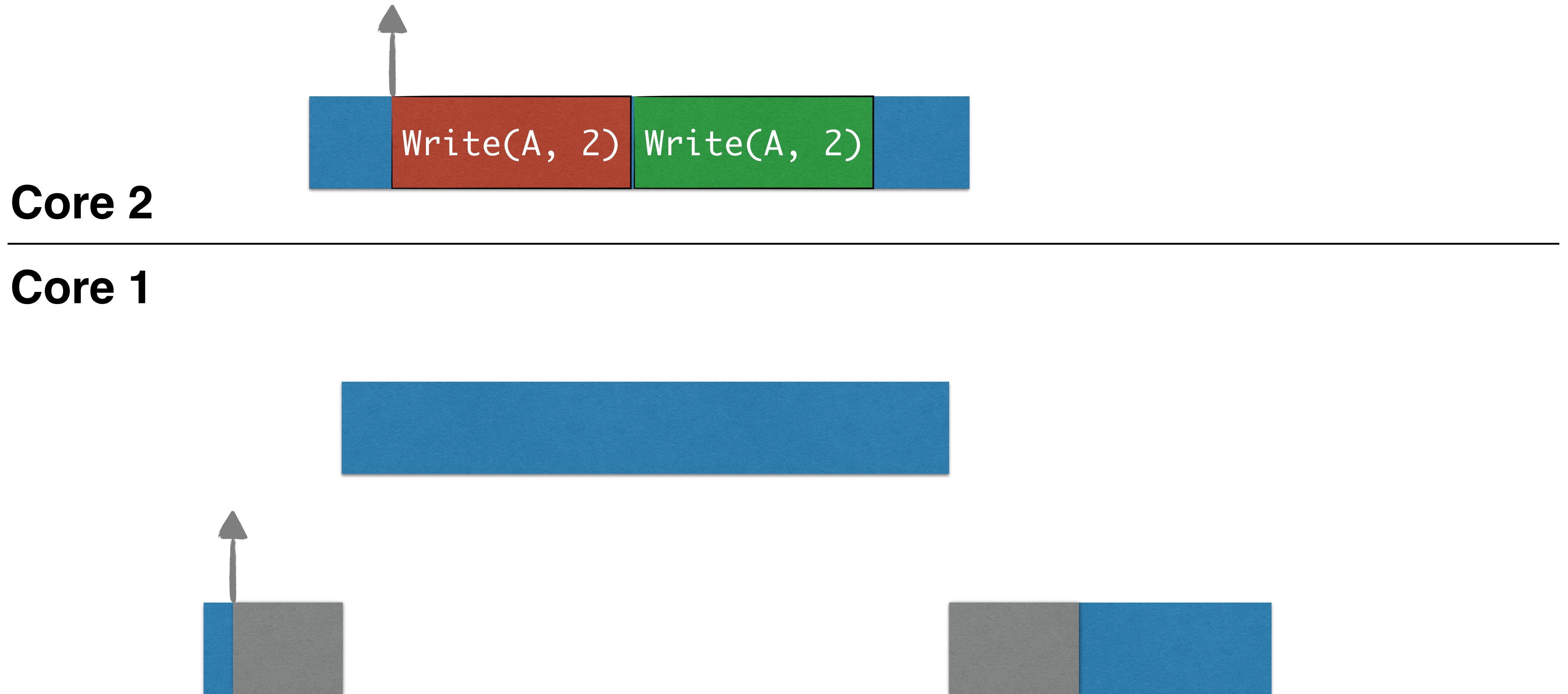
# Preemptions and serialisation



**Core 2**
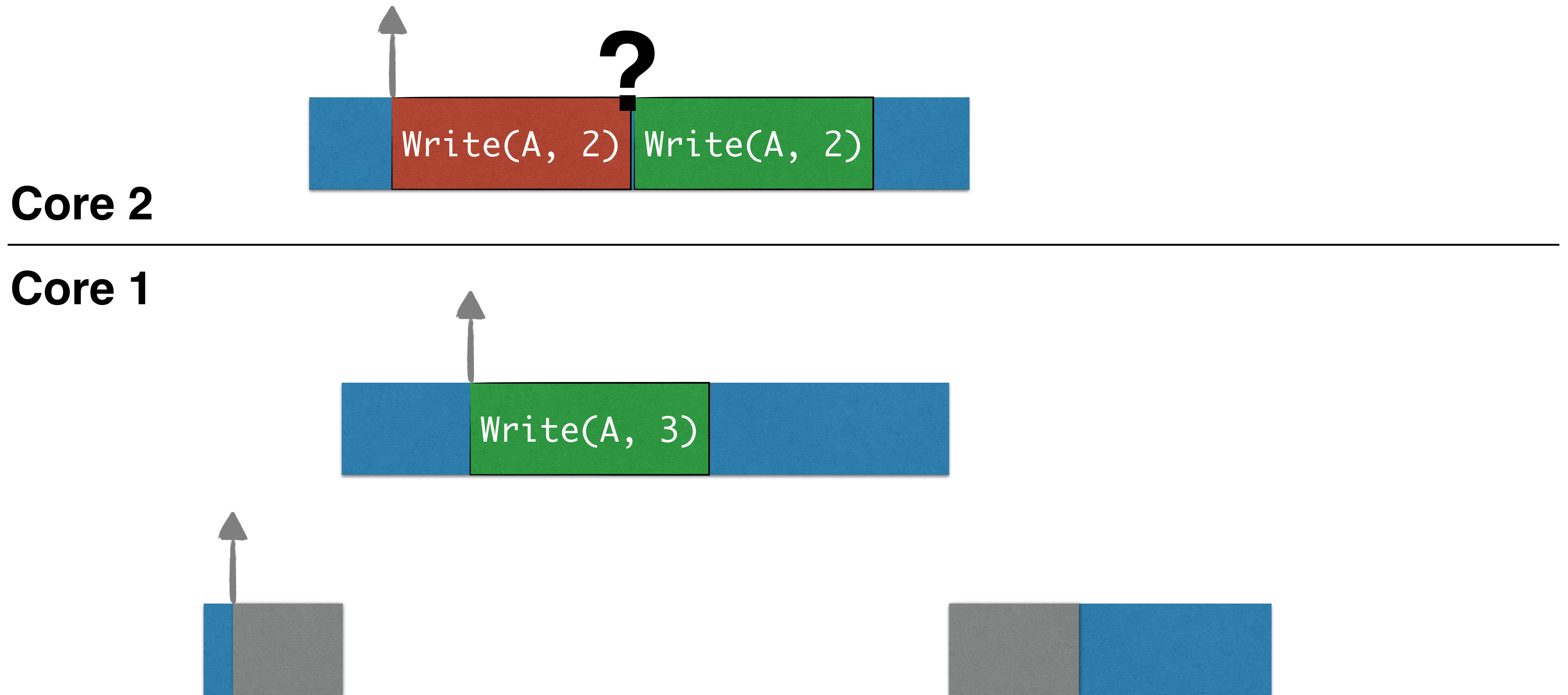
**Core 1**

# Preemptions and serialisation

Write(A, 2) Write(A, 2)

**Core 2**

**Core 1**

# Preemptions and serialisation



**Core 2**

**Core 1**

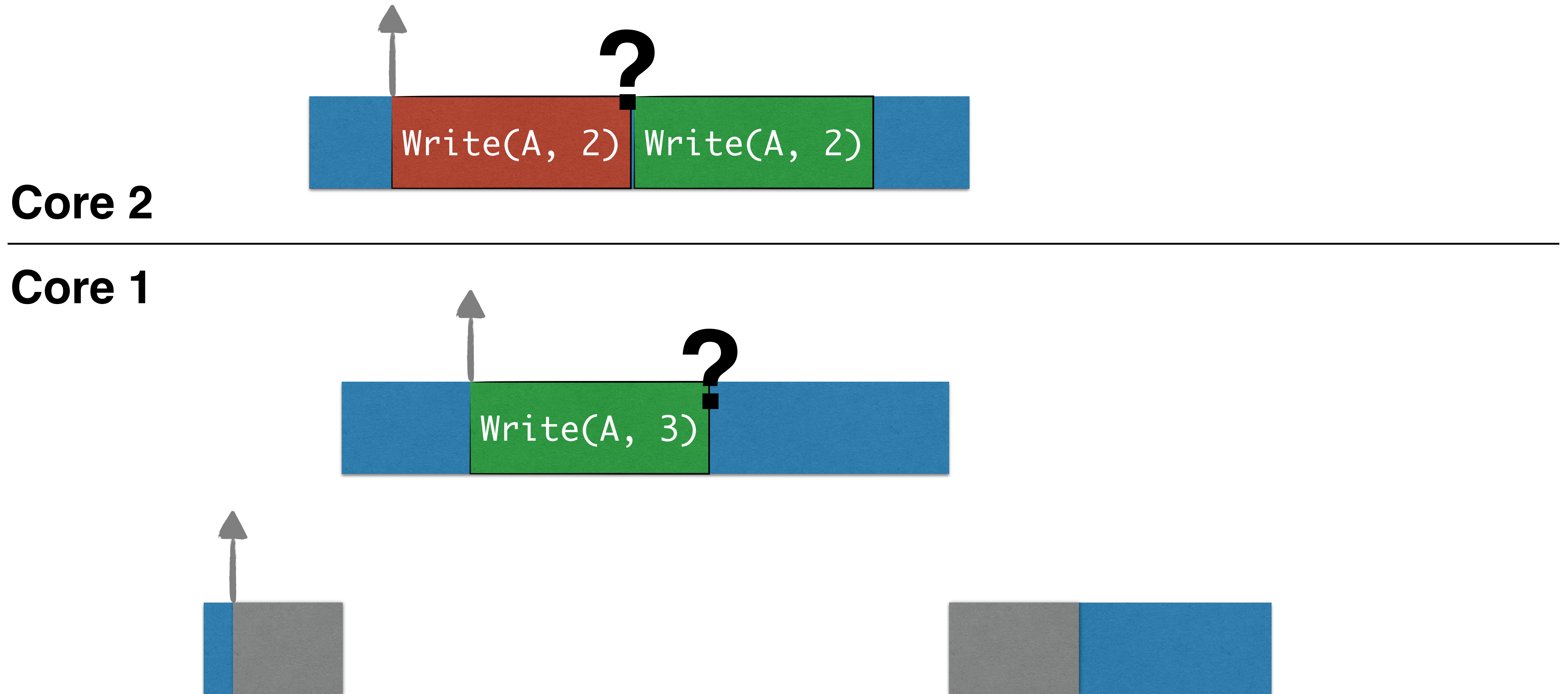# Preemptions and serialisation

# Preemptions and serialisation



**Core 2**

**Core 1**

# What to do?

Increase resistance to preemptions if a transaction can affect concurrent parallel transactions in jobs, while meeting all timing requirements.

Restrict to, at most, ONE transaction in progress, per core.

- No deadlocks.

- No transgression to FIFO serialisation.

# What to do?

Increase resistance to preemptions if a transaction can affect concurrent parallel transactions in jobs, while meeting all timing requirements.

Restrict to, at most, ONE transaction in progress, per core.

- No deadlocks.

- No transgression to FIFO serialisation.

SRP-TM

# Scheduling jobs with transactions: SRP-TM

# Assumptions

General scheduling rule: P-EDF.

While a transaction is in progress on a core: SRP ⇒ SRP-TM.

- Adds static preemption levels to tasks.

- Adds static preemption level to transactions.

- Adds variable ceiling to cores.

  - Highest preemption level of a task that could be waiting for the current transaction in progress to commit.

# Assigning preemption levels to tasks

Just like SRP, assign preemption levels to all tasks in set by increasing order of relative deadline…

… independently of core affinities.

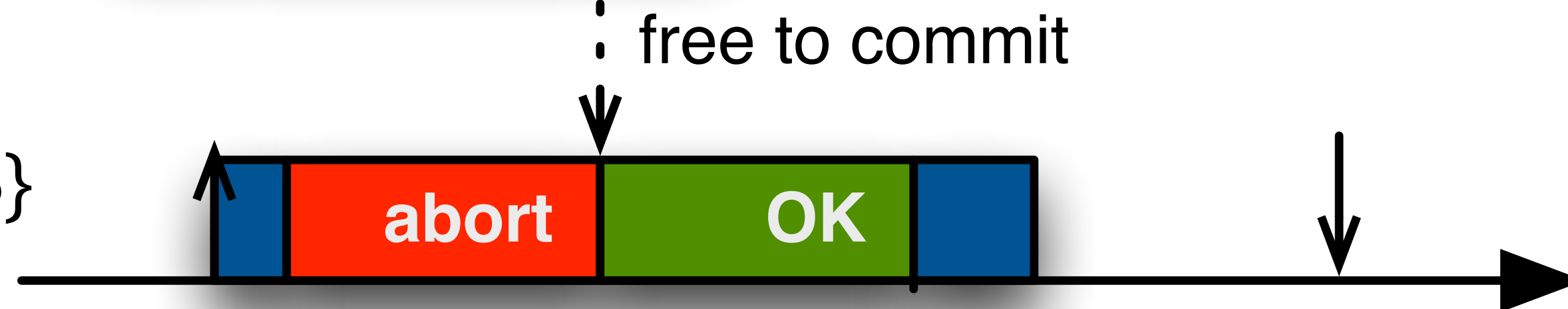| Task | Relative deadline | Preemption level |
|---|---|---|
| T5 | 120 | 1 |
| T2 | 100 | 2 |
| T3 | 80 | 3 |
| T4 | 70 | 4 |
| T6, T7 | 60 | 5 |
| T1 | 50 | 6 |

# Assigning preemption levels to transaction

Assign to each transaction the highest preemption level from all tasks that have one transaction that may depend on it to progress.
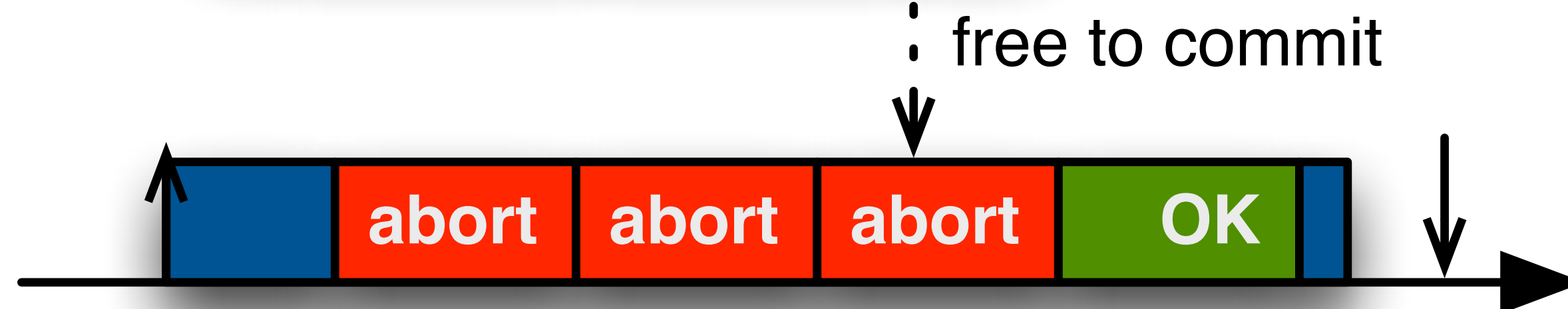
$\tau_1$ @ core 1
DS$_1$ = {A}

**OK**

free to commit

$\tau_2$ @ core 2
DS$_2$ = {A, B}

**abort**  **OK**

free to commit

$\tau_3$ @ core 3
DS$_3$ = {B}

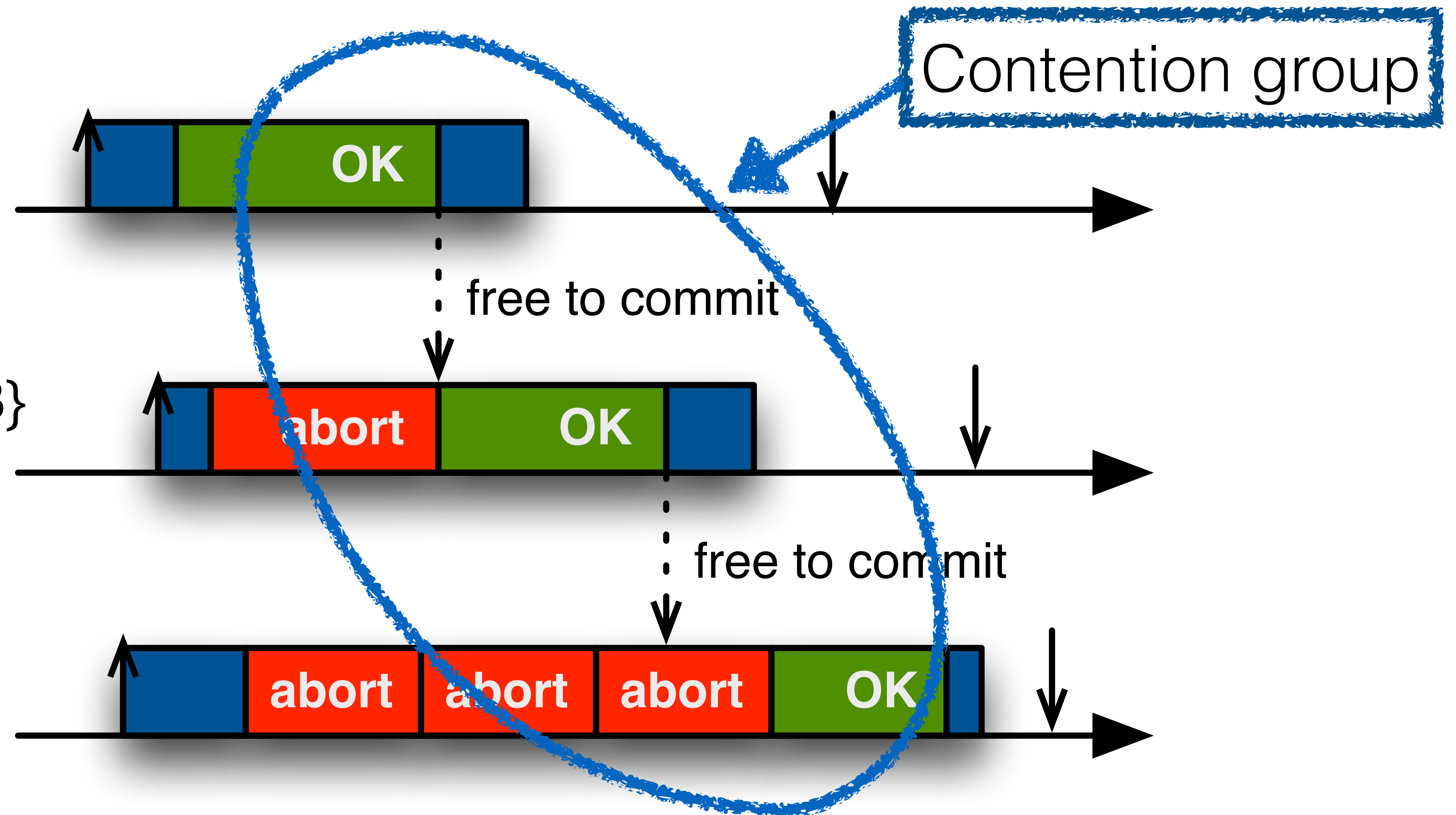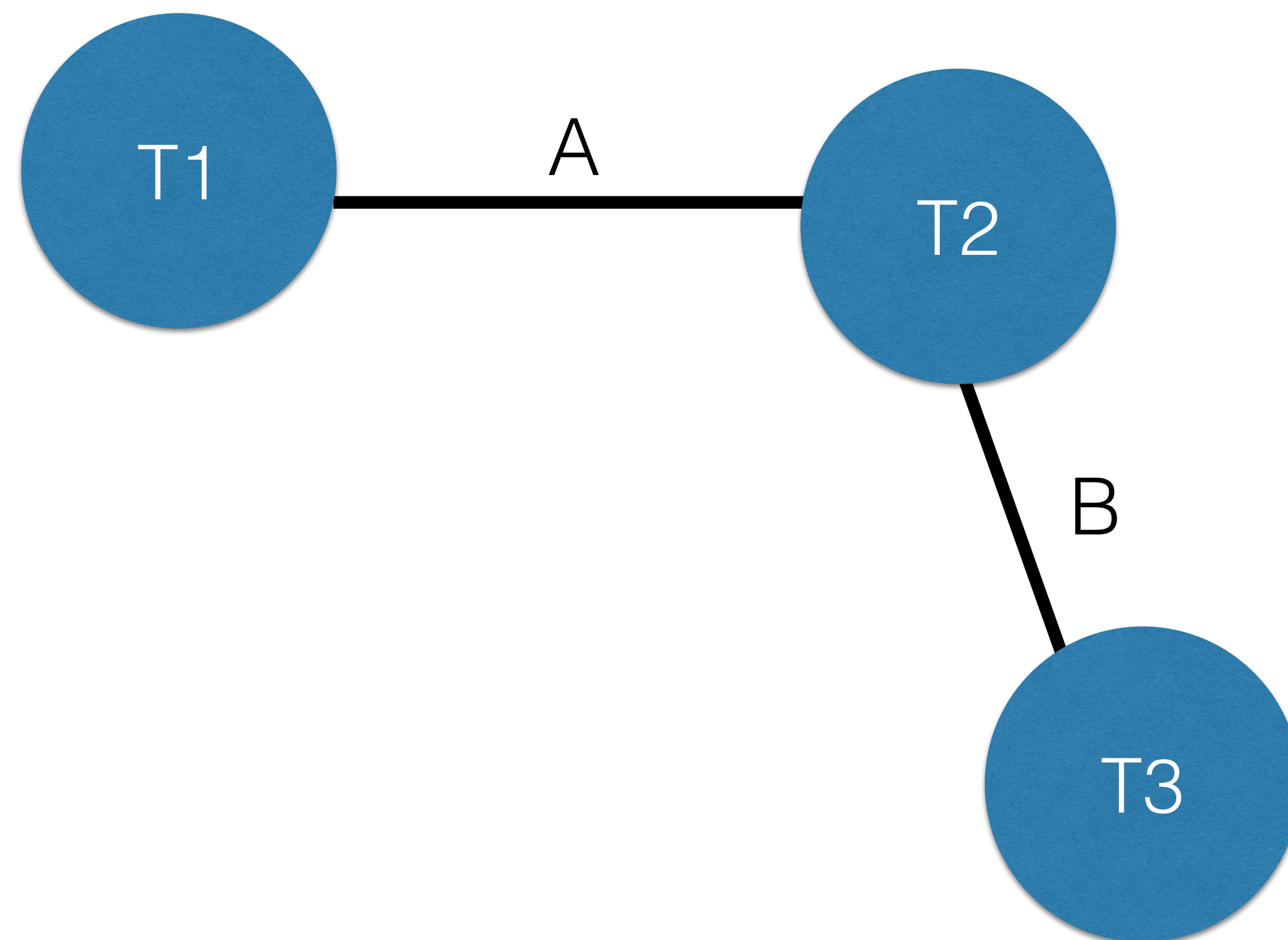**abort**  **abort**  **abort**  **OK**

# Assigning preemption levels to transaction

Assign to each transaction the highest preemption level from all tasks that have one transaction that may depend on it to progress.

# Assigning preemption levels to transaction

Assign to each transaction the highest preemption level from all tasks that have one transaction that may depend on it to progress.
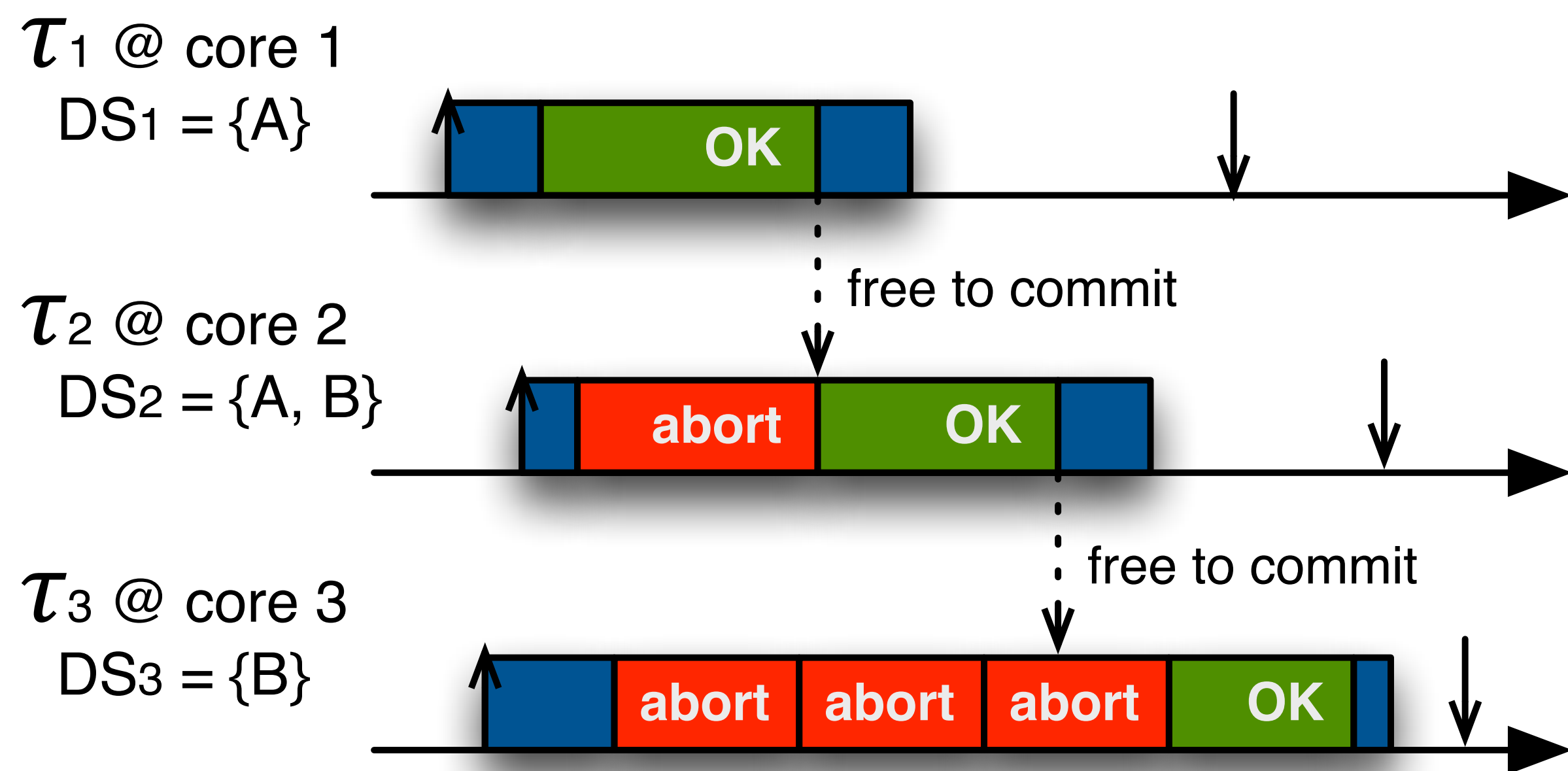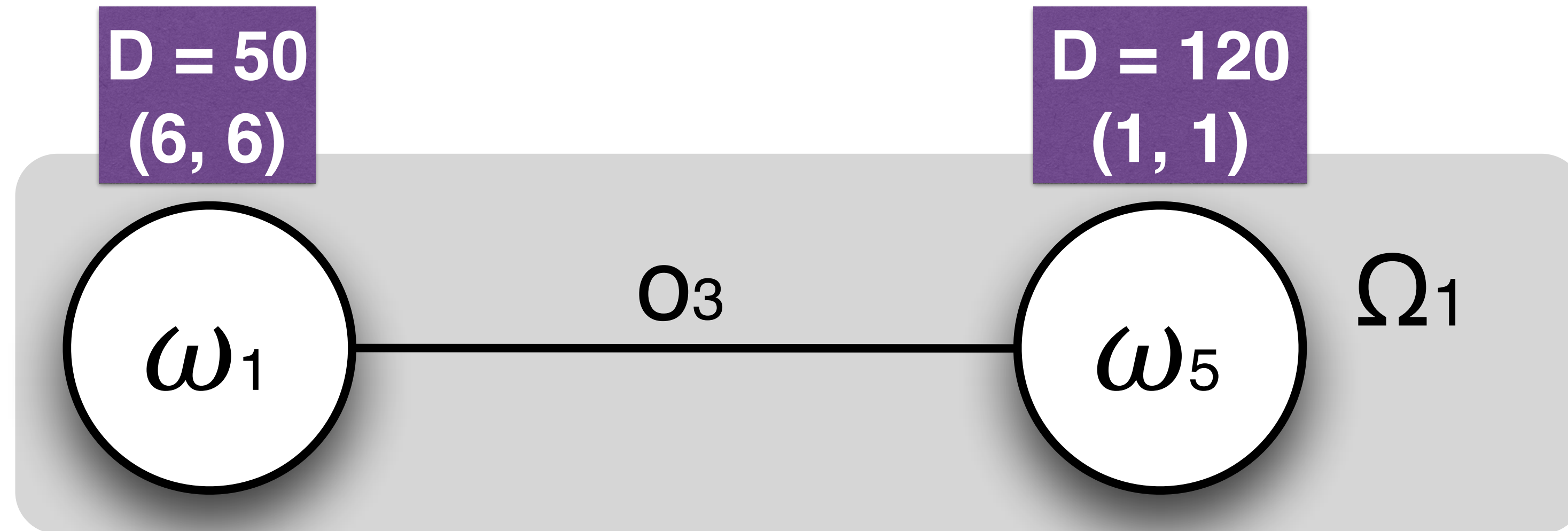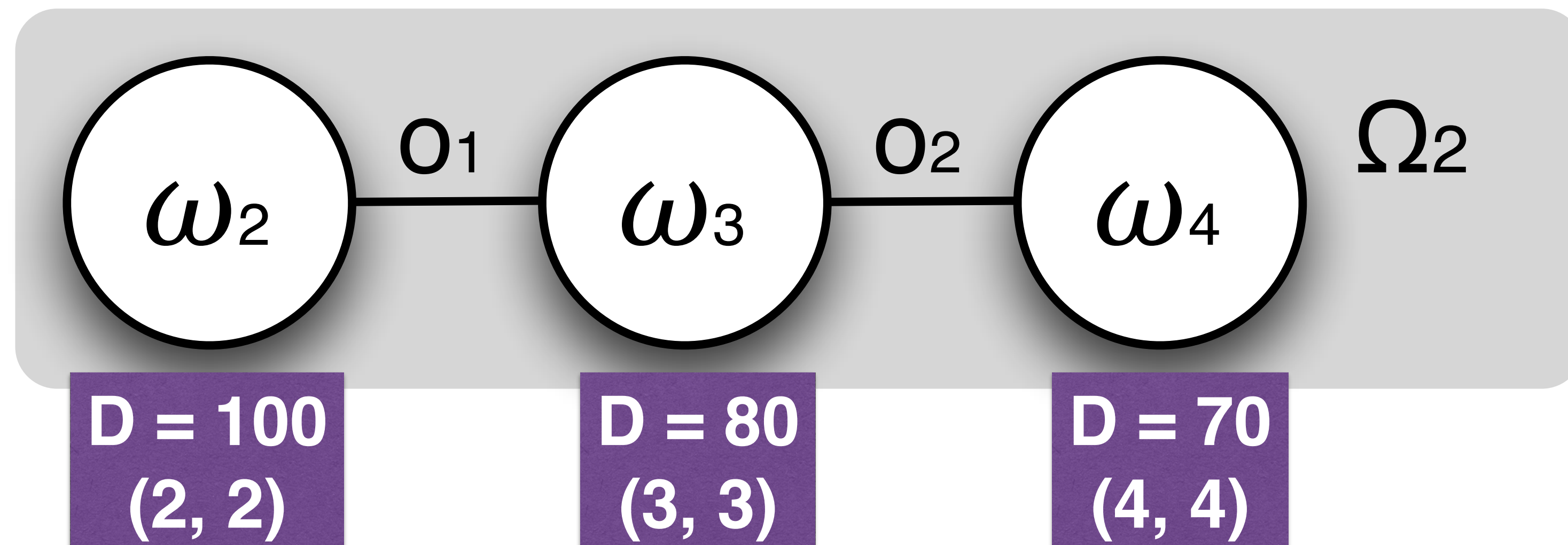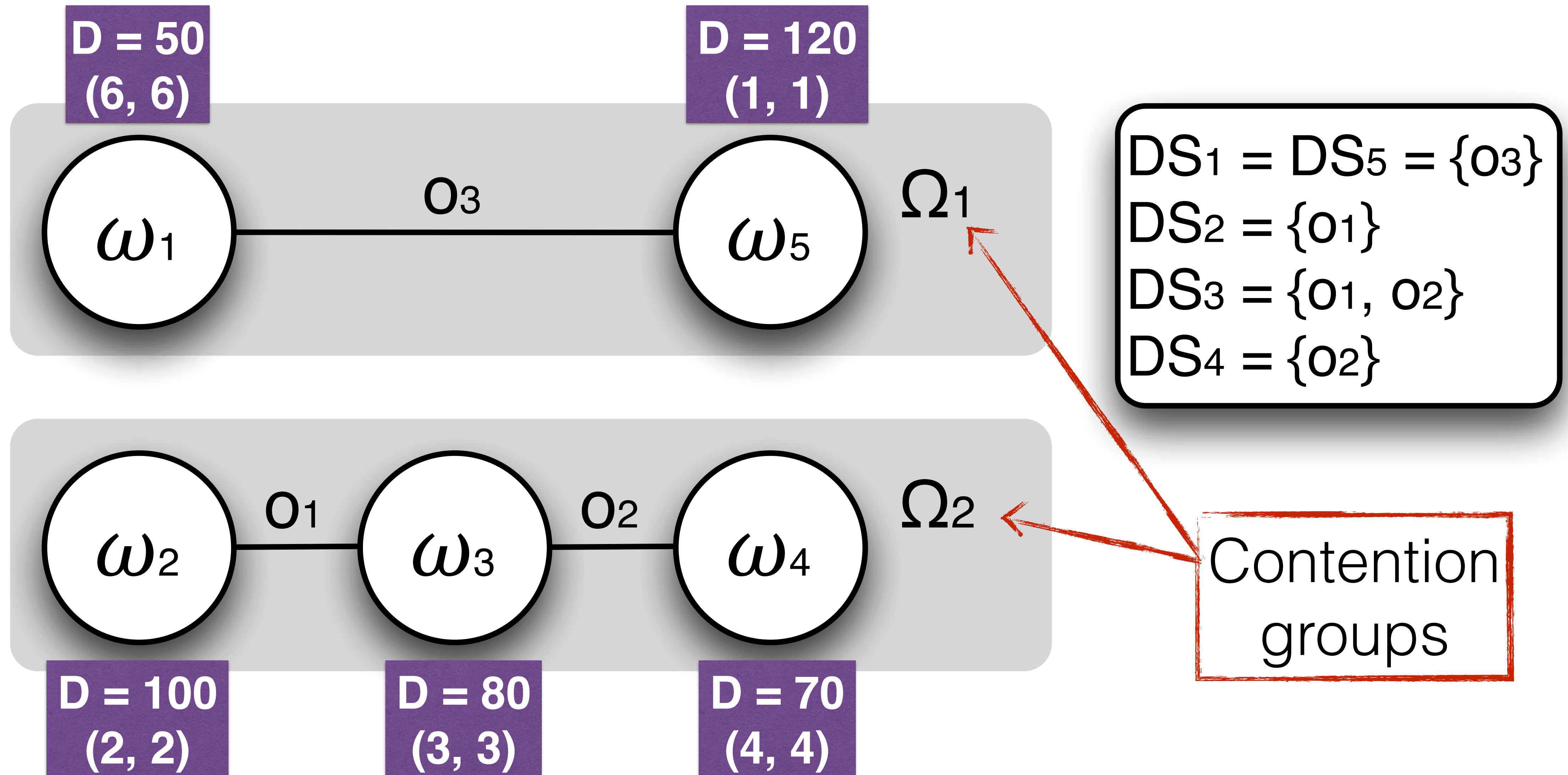
# A practical example

# A practical example

# A practical example

# Transaction vs. Transactionless

**P-EDF**

# Transaction vs. Transactionless

**P-EDF**



Transaction

# Transaction vs. Transactionless

**P-EDF**

**SRP-TM**

(8, 0)

(1, 6)

| | | |
|---|---|---|
| 0 | 6 | 0 |

(Task PL, Transaction PL)

X  Core ceiling

# Transaction vs. Transactionless

**P-EDF**

**SRP-TM**

(8, 0)

(1, 6)

| 0 | 6 | 0 |

**SRP-TM**

(4, 0)

(1, 6)

| 0 | 6 | 0 |

(Task PL, Transaction PL)

| X | Core ceiling |

# Transaction vs. Transactionless

**P-EDF**

**SRP-TM**

(8, 0)

(1, 6)

| 0 | 6 | 0 |

Direct blocking

**SRP-TM**

(4, 0)

(1, 6)

(Task PL, Transaction PL)

| X | Core ceiling |

| 0 | 6 | 0 |

# Transaction vs. Transaction

**P-EDF**

# Transaction vs. Transaction

**P-EDF**



**SRP-TM**

$(8, 0)$

$(1, 6)$

| 0 | 6 | 8 | 0 |

# Transaction vs. Transaction

**P-EDF**

**SRP-TM**

Direct blocking

(8, 0)

(1, 6)

| 0 | 6 | 8 | 0 |

# Mixing all together

**P-EDF**

# Mixing all together

**P-EDF**

**SRP-TM**

(8, 8)

(7, 0)

(1, 6)

| 0 | 6 | 8 | 0 | 8 | 0 |

# Mixing all together

**P-EDF**

Direct blocking

Indirect blocking

**SRP-TM**

(8, 8)

(7, 0)

(1, 6)

| 0 | 6 | 8 | 0 | 8 | 0 |
|---|---|---|---|---|---|

# SRP-TM operations in short

Transaction starts:

- Core ceiling is set to the preemption level of the transaction.

Transaction commits:

- Core ceiling is reset to zero.

# SRP-TM scheduling decisions in short

- Job in front of ready queue has transaction:

  - Core ceiling is raised to the preemption level of this task.

  - Job with transaction in progress is executed on behalf of job in front of ready queue.

- Job in front of ready queue does not have transaction:

  - Preempt running job iff has **earlier absolute deadline** than running job, and **higher preemption level** than core ceiling.

# Response time of a **transaction**

# Response time of a transaction



$\tau_1$ @ core 1
DS$_1$ = {A}

OK

free to commit

$\tau_2$ @ core 2
DS$_2$ = {A, B}

abort  OK

free to commit

$\tau_3$ @ core 3
DS$_3$ = {B}

abort  abort  abort  OK

# Response time of a transaction

$\tau_1$ @ core 1
DS$_1$ = {A}



free to commit

$\tau_2$ @ core 2
DS$_2$ = {A, B}

OK

abort OK

free to commit

$\tau_3$ @ core 3
DS$_3$ = {B}

abort abort abort OK

Transaction response time…

# Response time of a transaction

$\tau_1$ @ core 1
DS$_1$ = {A}

**OK**

free to commit

$\tau_2$ @ core 2
DS$_2$ = {A, B}

**abort**  **OK**

free to commit

$\tau_3$ @ core 3
DS$_3$ = {B}

**abort**  **abort**  **abort**  **OK**

Transaction response time…

… depends on
parallel transactions

# Response time of a transaction



$\tau_1$ @ core 1
DS$_1$ = {A}

OK

free to commit

$\tau_2$ @ core 2
DS$_2$ = {A, B}

abort    OK

free to commit

$\tau_3$ @ core 3
DS$_3$ = {B}

abort    abort    abort    OK

Transaction response time…

… depends on parallel transactions

… depends on intra-core interference

# Response time of a transaction

$\tau_1$ @ core 1
DS$_1$ = {A}

$\tau_2$ @ core 2
DS$_2$ = {A, B}

$\tau_3$ @ core 3
DS$_3$ = {B}

**OK**

**abort** **OK**

**abort** **abort** **abort** **OK**

free to commit

free to commit

Once it is free to commit, 2 more attempts, at most.

Transaction response time…

… depends on parallel transactions

… depends on intra-core interference

# Response time of a transaction

- The response time of the last transaction in a sequence of transactions is upper bounded by the sum of the response time of the last two attempts, for each transaction in the sequence.

# Response time of a transaction

- The response time of the last transaction in a sequence of transactions is upper bounded by the sum of the response time of the last two attempts, for each transaction in the sequence.

# Response time of a transaction

- The response time of the last two transactions depends exclusively on intra-core interference:

  - **IT CAN BE ANALYTICALLY UPPER BOUNDED!**

- Maximum response time of a transaction…

  - Determine every possible sequence, sum response times of last two attempts and choose the maximum value… **COMBINATIONAL ORDER!!!**

  - For every processor, choose the maximum response time of last two attempts of a transaction that belongs to the same contention group, and sum them all… **PESSIMISTIC, but LINEAR ORDER!**

# Response time of a **task**

# Blocking and interference

# Blocking and interference

# Blocking and interference



Direct Blocking + Interference

Interference

Interference of jobs without transactions

Interference

Transaction

$r_{i,j}$

$d_{i,j}$

Interference + IB

DB + I

Interference

$r_{i,j}$

$d_{i,j}$

IB:  Indirect blocking
DB: Direct blocking
I:    Interference

# Simulation results

# Simulation conditions

- Scheduling policies:

  - pure P-EDF

  - NPUC

  - NPDA

  - SRP-TM

  - FLMP

# Simulation conditions

- Experiment 1: varying system size

  - Variable number of cores: $m \in \{2, 4, 8, 16, 32, 64\}$

  - Number of transactional objects linear with $m$: $p \in \{5, 10, 20, 40, 80, 160\}$, so each object is accessed by 3 task, on average.

  - Number of contention groups linear with $m$: $g \in \{1, 2, 4, 8, 16, 32\}$, so each group maintains the same size and the same expected number of tasks.

# Simulation conditions

- Experiment 2: varying size of contention groups

  - Constant number of cores: $m = 64$.

  - Constant number of transactional objects linear: $p = 160$.

  - Variable number of contention groups: $g \in \{1, 2, 4, 8, 16, 32\}$, so to observe the effects of granularity of contention groups for systems with same size.

# Feasibility (experiment 1)



Legend:
- PEDF
- NPUC
- NPDA
- SRPTM
- FMLP

X-axis: (2, 1), (4, 2), (8, 4), (16, 8), (32, 16), (64, 32)

(Number of cores, number of groups)

Y-axis: 0%, 10%, 20%, 30%, 40%

| | | | | | |
|---|---|---|---|---|---|
| **(8, 4)** | 80% | 68% | 86% | 82% | 38% |
| **(16, 8)** | 66% | 30% | 68% | 60% | 6% |
| **(32, 16)** | 34% | 0% | 12% | 12% | 0% |
| **(64, 32)** | 2% | 0% | 0% | 0% | 0% |

| | | | | | |
|---|---|---|---|---|---|
| **(64, 16)** | 4% | 0% | 0% | 0% | 0% |
| **(64, 8)** | 0% | 0% | 0% | 0% | 0% |
| **(64, 4)** | 6% | 0% | 0% | 2% | 0% |
| **(64, 2)** | 6% | 0% | 6% | 2% | 0% |
| **(64, 1)** | 4% | 0% | 2% | 2% | 0% |

it
e



Experiment 1



Experiment 2

| | | | | | |
|---|---|---|---|---|---|
| **(8, 4)** | 34826 | 35215 | 34696 | 30713 | 81666 |
| **(16, 8)** | 85609 | 95317 | 90535 | 68950 | 220041 |
| **(32, 16)** | 188286 | 202247 | 193040 | 156309 | 459582 |
| **(64, 32)** | 432176 | 438299 | 423635 | 373730 | 1006744 |

| | | | | | |
|---|---|---|---|---|---|
| **(64, 16)** | 443036 | 429845 | 413484 | 387422 | 1232219 |
| **(64, 8)** | 473085 | 454010 | 433336 | 404777 | 1344484 |
| **(64, 4)** | 475156 | 459013 | 439262 | 405838 | 1369213 |
| **(64, 2)** | 484011 | 444793 | 430496 | 421342 | 1341424 |
| **(64, 1)** | 503295 | 471157 | 454371 | 438009 | 1383684 |



Experiment 1



Experiment 2

| | | | | | |
|---|---|---|---|---|---|
| **(8, 4)** | 1221 | 2562 | 1948 | 607 | 5718 |
| **(16, 8)** | 3225 | 6973 | 6549 | 1535 | 17392 |
| **(32, 16)** | 6425 | 17433 | 14948 | 5258 | 45713 |
| **(64, 32)** | 28501 | 53829 | 49013 | 23435 | 179048 |

| | | | | | |
|---|---|---|---|---|---|
| **(64, 16)** | 24273 | 51199 | 46335 | 17629 | 622868 |
| **(64, 8)** | 33908 | 55863 | 48934 | 24574 | 1200480 |
| **(64, 4)** | 26135 | 55589 | 51026 | 17298 | 1346441 |
| **(64, 2)** | 25766 | 48684 | 43582 | 17525 | 1338480 |
| **(64, 1)** | 26879 | 57208 | 49244 | 20375 | 1382911 |



Experiment 1



Experiment 2

| | | | | | |
|---|---|---|---|---|---|
| **(4, 2)** | 24% | 21% | 21% | 13% | 33% |
| **(8, 4)** | 29% | 31% | 31% | 20% | 40% |
| **(16, 8)** | 32% | 58% | 59% | 21% | 63% |
| **(32, 16)** | 32% | 58% | 58% | 22% | 63% |
| **(64, 32)** | 33% | 62% | 62% | 23% | 69% |

| | | | | | |
|---|---|---|---|---|---|
| **(64, 16)** | 35% | 65% | 66% | 25% | 141% |
| **(64, 8)** | 36% | 67% | 67% | 26% | 323% |
| **(64, 4)** | 36% | 68% | 69% | 26% | 682% |
| **(64, 2)** | 37% | 67% | 67% | 27% | 1 451% |
| **(64, 1)** | 37% | 70% | 70% | 26% | 3 074% |



Experiment 1



Experiment 2

| | | | | | |
|---|---|---|---|---|---|
| (4, 2) | 24% | 21% | 21% | 13% | 33% |
| (8, 4) | 29% | 31% | 31% | 20% | 40% |
| (16, 8) | 32% | 58% | 59% | 21% | 63% |
| (32, 16) | 32% | 58% | 58% | 22% | 63% |
| (64, 32) | 33% | 62% | 62% | 23% | 69% |

| | | | | | |
|---|---|---|---|---|---|
| (64, 8) | 36% | 67% | 67% | 26% | 323% |
| (64, 4) | 36% | 68% | 69% | 26% | 682% |
| (64, 2) | 37% | 67% | 67% | 27% | 1 451% |
| (64, 1) | 37% | 70% | 70% | 26% | 3 074% |



Experiment 1



Experiment 2

# Wrapping up

# Conclusion (1/2)

- FIFO serialisation is the predictable and fair.

- Scheduling has an effect on the performance of transactions.

  - SRP-TM extends P-EDF when a transaction is in progress.

    - Takes into account **possible** concurrent parallel transactions with earlier deadlines, **without sharing** scheduling data between cores.

    - Allows jobs with earlier deadlines to **preempt** or **speed up** a transaction in progress.

# Conclusion (2/2)

- We provide an analytical method to upper bound the response time of transactions under SRP-TM.

- We provide an analytical method to upper bound the response time of tasks under SRP-TM.

# That's it! Thanks! Questions?