# Resource Management for Mixed-Criticality Systems on Multi-Core Platforms with Focus on Communication

## Embedded Tutorial

Robin Arbaud, Dávid Juhász, Axel Jantsch

DSD 2018, Prague, Czech Republic

August 31, 2018
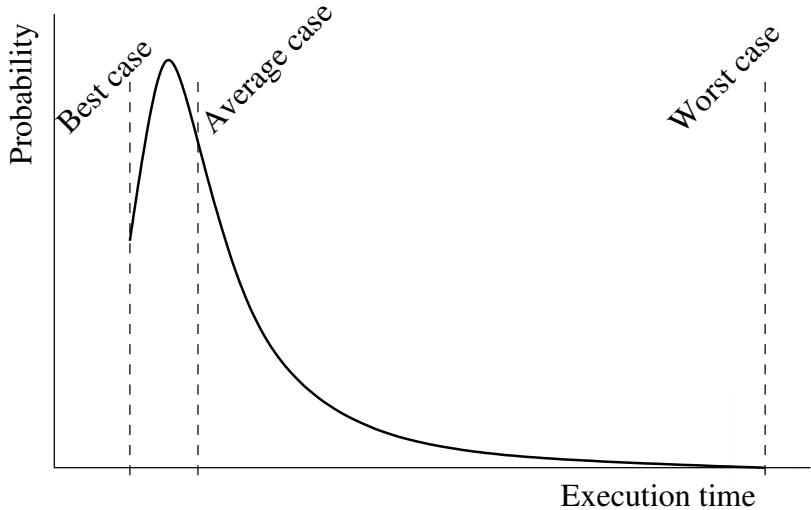
# Outline

www.ict.tuwien.ac.at

Types of tasks:

Best Effort  May take as long as needed;

Soft Real-Time  Have deadlines, but may miss some deadlines;

Hard Real-time  No deadline must be missed, ever.

- **Full Isolation** to accommodate the worst case:
  - Dedicated resources for a given task;
  - Addresses real-time and safety requirements;
  - Costly due to over-provisioning;
  - Complete isolation = no interference.

- **Full Isolation** to accommodate the worst case:
  - Dedicated resources for a given task;
  - Addresses real-time and safety requirements;
  - Costly due to over-provisioning;
  - Complete isolation = no interference.
- **Full sharing** to accommodate the average case:
  - Best average performance for given cost;
  - Lowest cost for target average performance;
  - Highest efficiency;
  - Unbounded worst case performance.
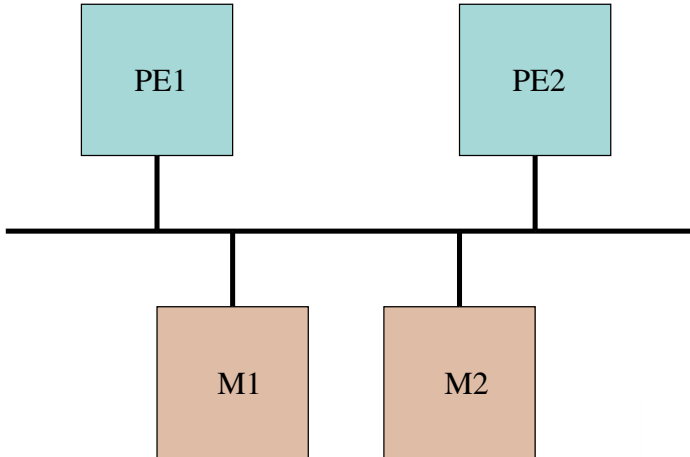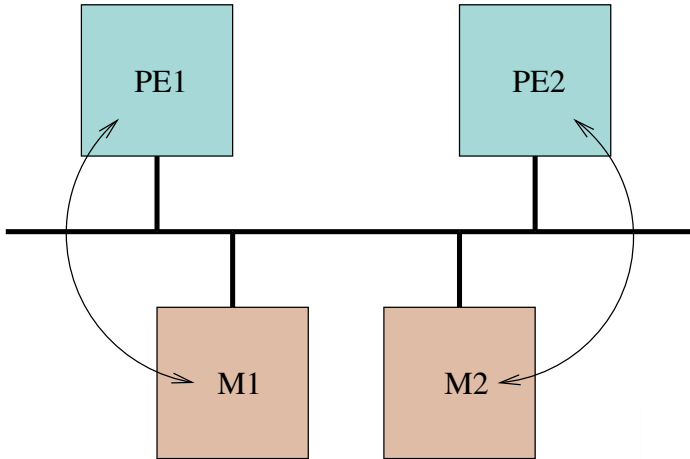
www.ict.tuwien.ac.at

- **Full Isolation** to accommodate the worst case:
  - Dedicated resources for a given task;
  - Addresses real-time and safety requirements;
  - Costly due to over-provisioning;
  - Complete isolation = no interference.
- **Full sharing** to accommodate the average case:
  - Best average performance for given cost;
  - Lowest cost for target average performance;
  - Highest efficiency;
  - Unbounded worst case performance.
- Sharing with **bounded interference**:
  - Allowing and controlling interference can provide real-time and safety requirement;
  - Reduced costs;
  - Applicable for mix of critical and best effort tasks.

www.ict.tuwien.ac.at

- Direct interference and competition for the same resource;
- Indirect interference through shared resources:
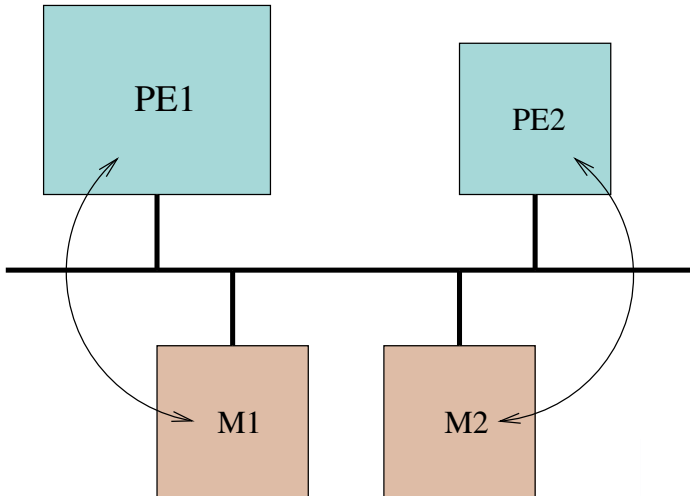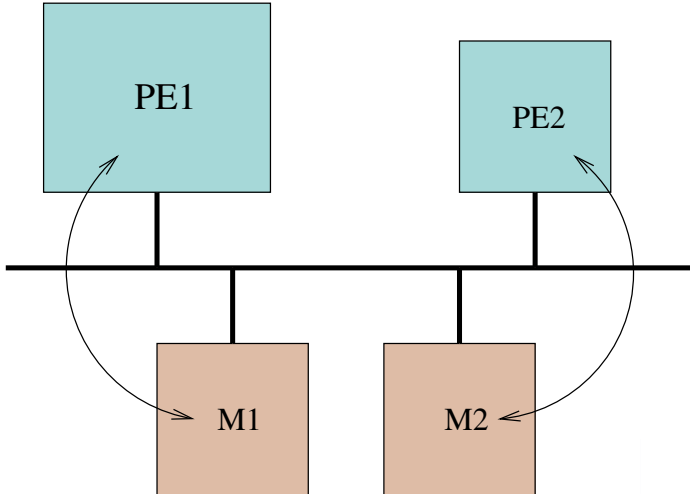  - Performance inversion,
  - Over-synchronization.

- PE1 communicates with M1 and PE2 communicates with M2;
- All tasks meet deadlines;

- Replacing PE1 with a faster PE;

- Replacing PE1 with a faster PE;
- **May increase execution time of PE2 tasks.**

# Over-Synchronization



- Assumption: Bounded buffers between tasks;
- No control or data dependency between C and D branches;

- Assumption: Bounded buffers between tasks;
- No control or data dependency between C and D branches;

- Assumption: Bounded buffers between tasks;
- No control or data dependency between C and D branches;

- Assumption: Bounded buffers between tasks;
- No control or data dependency between C and D branches;
- **If C is delayed or stuck, D suffers.**

- Spatial isolation: No sharing at any time;
- Temporal isolation: Sharing at pre-defined time periods.

# Temporal Isolation

- Many resources lead to a huge design space;

- Many resources lead to a huge design space;
- Types of Resources:

- Many resources lead to a huge design space;
- Types of Resources:
  - Computation: Processing elements,

- Many resources lead to a huge design space;
- Types of Resources:
  - Computation: Processing elements,
  - Storage: Buffers, caches, off-chip memory,

- Many resources lead to a huge design space;
- Types of Resources:
  - Computation: Processing elements,
  - Storage: Buffers, caches, off-chip memory,
  - **Communication: Buses, links, NoCs.**

www.ict.tuwien.ac.at

Real-Time Constraints from Event to System Response

www.ict.tuwien.ac.at

## Real-Time Constraints from Event to System Response

- Deadlines based on real-time requirements

## Real-Time Constraints from Event to System Response

- Deadlines based on real-time requirements
- No deadline misses allowed ever

### Real-Time Constraints from Event to System Response

- Deadlines based on real-time requirements
- No deadline misses allowed ever
- Prepare for the worst-case

## Task Set

| Task | Period | Deadline | WCET |
|------|--------|----------|------|
| $T1$ | 10 | 5 | 4 |
| $T2$ | 10 | 10 | 6 |

## Task Set

| Task | Period | Deadline | WCET |
|------|--------|----------|------|
| $T1$ | 10     | 5        | 4    |
| $T2$ | 10     | 10       | 6    |



time

## Task Set

| Task | Period | Deadline | WCET | ACET |
|------|--------|----------|------|------|
| $T1$ | 10 | 5 | 4 | 2 |
| $T2$ | 10 | 10 | 6 | 4 |

## Task Set

| Task | Period | Deadline | WCET | ACET |
|------|--------|----------|------|------|
| $T1$ | 10 | 5 | 4 | 2 |
| $T2$ | 10 | 10 | 6 | 4 |



time

- Static guarantees even for the worst-case



time

- Static guarantees even for the worst-case
- Overprovisioning for the average-case

- Static guarantees even for the worst-case
- Overprovisioning for the average-case



time

**Economic need for trading off guarantees for utilization**

www.ict.tuwien.ac.at

Worst-Case Execution Time Estimates

www.ict.tuwien.ac.at

**Worst-Case Execution Time Estimates**

- Various components are validated against various assumptions

## Worst-Case Execution Time Estimates

- Various components are validated against various assumptions
- The stricter assumptions, the longer estimated WCET

## Worst-Case Execution Time Estimates

- Various components are validated against various assumptions
- The stricter assumptions, the longer estimated WCET

## Criticality Levels of Tasks

## Worst-Case Execution Time Estimates

- Various components are validated against various assumptions
- The stricter assumptions, the longer estimated WCET

## Criticality Levels of Tasks

- Strictness of assumptions

## Worst-Case Execution Time Estimates

- Various components are validated against various assumptions
- The stricter assumptions, the longer estimated WCET

## Criticality Levels of Tasks

- Strictness of assumptions
- Ordered set of levels

- Start system in the lowest criticality mode

- Start system in the lowest criticality mode
- Schedule tasks with criticality not below the current criticality mode

- Start system in the lowest criticality mode
- Schedule tasks with criticality not below the current criticality mode
- Monitor whether execution violates current assumptions

- Start system in the lowest criticality mode
- Schedule tasks with criticality not below the current criticality mode
- Monitor whether execution violates current assumptions
- Switch to higher criticality mode for stricter assumptions

## Task Set *HI*

| Task | Period | Deadline | WCET |
|------|--------|----------|------|
| $T1$ | 10     | 5        | 4    |
| $T2$ | 10     | 10       | 6    |

## Task Set *HI*

| Task | Period | Deadline | WCET |
|------|--------|----------|------|
| T1   | 10     | 5        | 4    |
| T2   | 10     | 10       | 6    |

## Task Set *LO*

| Task | Period | Deadline | WCET |
|------|--------|----------|------|
| T1   | 10     | 5        | 2    |
| T2   | 10     | 10       | 4    |
| T3   | 10     | 9        | 3    |

## Mixed-Criticality Task Set

| Task | CL | Period | Deadline | WCET |
|------|-----|--------|----------|--------|
| $T1$ | $HI$ | 10 | 5 | $[2, 4]$ |
| $T2$ | $HI$ | 10 | 10 | $[4, 6]$ |
| $T3$ | $LO$ | 10 | 9 | $[3]$ |

## Mixed-Criticality Task Set

| Task | CL | Period | Deadline | WCET |
|------|------|--------|----------|--------|
| $T1$ | $HI$ | 10 | 5 | $[2,4]$ |
| $T2$ | $HI$ | 10 | 10 | $[4,6]$ |
| $T3$ | $LO$ | 10 | 9 | $[3]$ |



time

## Mixed-Criticality Task Set

| Task | CL | Period | Deadline | WCET |
|------|----|--------|----------|------|
| $T1$ | $HI$ | 10 | 5 | $[2, 4]$ |
| $T2$ | $HI$ | 10 | 10 | $[4, 6]$ |
| $T3$ | $LO$ | 10 | 9 | $[3]$ |

www.ict.tuwien.ac.at

## Mixed-Criticality Task Set

| Task | CL | Period | Deadline | WCET |
|------|-----|--------|----------|--------|
| $T1$ | $HI$ | 10 | 5 | $[2, 4]$ |
| $T2$ | $HI$ | 10 | 10 | $[4, 6]$ |
| $T3$ | $LO$ | 10 | 9 | $[3]$ |



**Ignoring *LO* criticality tasks in *HI* mode**

## Mixed-Criticality Task Set

| Task | CL | Period | Deadline | WCET |
|------|-----|--------|----------|--------|
| $T1$ | $HI$ | 10 | 5 | $[2, 4]$ |
| $T2$ | $HI$ | 10 | 10 | $[4, 6]$ |
| $T3$ | $LO$ | 10 | 9 | $[3]$ |

**Ignoring $LO$ criticality tasks in $HI$ mode**

Static Verification

Runtime Robustness

## Static Verification

- Timing guarantees for each mode

## Runtime Robustness

## Static Verification

- Timing guarantees for each mode
- **No link to requirements of safety standards**

## Runtime Robustness

## Static Verification

- Timing guarantees for each mode
- **No link to requirements of safety standards**

## Runtime Robustness

- Switching mode to restrict assumptions as necessary

## Static Verification

- Timing guarantees for each mode
- **No link to requirements of safety standards**

## Runtime Robustness

- Switching mode to restrict assumptions as necessary
- **Ignoring low-criticality tasks**

## Static Verification

- Timing guarantees for each mode
- **No link to requirements of safety standards**

## Runtime Robustness

- Switching mode to restrict assumptions as necessary
- **Ignoring low-criticality tasks**
- **Returning to low-criticality mode**

# Resource Management for Mixed-Criticality Systems

www.ict.tuwien.ac.at

Cache Misses

## Cache Misses

- Memory needs to be accessed

### Cache Misses

- Memory needs to be accessed
- Critical tasks may be penalized by interference

## Cache Misses

- Memory needs to be accessed
- Critical tasks may be penalized by interference

## Limited Interference for Critical Tasks

# Memory Accesses and Execution Time

## Cache Misses

- Memory needs to be accessed
- Critical tasks may be penalized by interference

## Limited Interference for Critical Tasks

- Reduce WCET of critical tasks by limiting interference

## Cache Misses

- Memory needs to be accessed
- Critical tasks may be penalized by interference

## Limited Interference for Critical Tasks

- Reduce WCET of critical tasks by limiting interference
- Budget of cache misses for non-critical tasks (cores)

## Cache Misses

- Memory needs to be accessed
- Critical tasks may be penalized by interference

## Limited Interference for Critical Tasks

- Reduce WCET of critical tasks by limiting interference
- Budget of cache misses for non-critical tasks (cores)
- Suspend execution of tasks with depleted budget

## Cache Misses

- Memory needs to be accessed
- Critical tasks may be penalized by interference

## Limited Interference for Critical Tasks

- Reduce WCET of critical tasks by limiting interference
- Budget of cache misses for non-critical tasks (cores)
- Suspend execution of tasks with depleted budget

**Can be generalized to any shared resource**

## Pros

## Cons

## Pros

- Low overhead on COTS

## Cons

## Pros

- Low overhead on COTS
  - Performance counters
  - Interprocessor interrupts

## Cons

# Memory Access Budgeting – Pros and Cons

## Pros

- Low overhead on COTS
    - Performance counters
    - Interprocessor interrupts

## Cons

- Limited to 2 levels, critical and non-critical

## Pros

- Low overhead on COTS
  - Performance counters
  - Interprocessor interrupts

## Cons

- Limited to 2 levels, critical and non-critical
- No general algorithm for deriving budgets

## Pros

- Low overhead on COTS
  - Performance counters
  - Interprocessor interrupts

## Cons

- Limited to 2 levels, critical and non-critical
- No general algorithm for deriving budgets
- Not scaling with the number of non-critical cores

www.ict.tuwien.ac.at

WCET and Observation Points for Critical Tasks

## WCET and Observation Points for Critical Tasks

- Remaining WCET for each observation point

## WCET and Observation Points for Critical Tasks

- Remaining WCET for each observation point
- Worst-case interference between observation points

## WCET and Observation Points for Critical Tasks

- Remaining WCET for each observation point
- Worst-case interference between observation points

## Safety Condition at each Observation Point

The critical task will not miss its deadline even if worst-case happens until the next observation point

Evaluate safety condition at each observation point:

Evaluate safety condition at each observation point:

code segment $k-$

code segment $k$

code segment $k+$

Evaluate safety condition at each observation point:

code segment $k-$
observation point $k$
code segment $k$
observation point $k+1$
code segment $k+$

Evaluate safety condition at each observation point:

| code segment $k-$ | | execution time |
| observation point $k$ | | |
| code segment $k$ | $+$ | $WCET_{interference}(k, k+1)$ |
| observation point $k+1$ | | |
| code segment $k+$ | $+$ | $WCET_{isolation}(k+1, end)$ |
| | $+$ | observation overhead |
| | $<$ | deadline |

Evaluate safety condition at each observation point:

| | |
|---|---|
| code segment $k-$ | execution time |
| observation point $k$ | |
| code segment $k$ | $+ \quad WCET_{interference}(k, k+1)$ |
| observation point $k+1$ | |
| code segment $k+$ | $+ \quad WCET_{isolation}(k+1, end)$ |
| | |
| | $+ \quad$ observation overhead |
| | |
| | $< \quad$ deadline |

Suspend non-critical tasks if safety condition does not hold.

## Pros

## Cons

## Pros

- Code instrumentation is easy to implement on COTS

## Cons

## Pros

- Code instrumentation is easy to implement on COTS
- Monitoring execution time limits pessimism

## Cons

## Pros

- Code instrumentation is easy to implement on COTS
- Monitoring execution time limits pessimism
- Evaluating condition at observation points limits pessimism

## Cons

## Pros

- Code instrumentation is easy to implement on COTS
- Monitoring execution time limits pessimism
- Evaluating condition at observation points limits pessimism

## Cons

- Limited to 2 levels, critical and non-critical

## Pros

- Code instrumentation is easy to implement on COTS
- Monitoring execution time limits pessimism
- Evaluating condition at observation points limits pessimism

## Cons

- Limited to 2 levels, critical and non-critical
- Large execution time overhead of monitoring

# Resource Management for Mixed-Criticality Systems

Workload Arrival Function (WAF) for each Task

## Workload Arrival Function (WAF) for each Task

- Task is activated by events

## Workload Arrival Function (WAF) for each Task

- Task is activated by events
- The events are mapped to WCET values

## Workload Arrival Function (WAF) for each Task

- Task is activated by events
- The events are mapped to WCET values
- WAF accumulates required execution time

## Workload Arrival Function (WAF) for each Task

- Task is activated by events
- The events are mapped to WCET values
- WAF accumulates required execution time

## Arrival Monitoring and Admission

## Workload Arrival Function (WAF) for each Task

- Task is activated by events
- The events are mapped to WCET values
- WAF accumulates required execution time

## Arrival Monitoring and Admission

- Actual workload is calculated based on WAFs

## Workload Arrival Function (WAF) for each Task

- Task is activated by events
- The events are mapped to WCET values
- WAF accumulates required execution time

## Arrival Monitoring and Admission

- Actual workload is calculated based on WAFs
- Check whether workload would exceed serviceable level

## Workload Arrival Function (WAF) for each Task

- Task is activated by events
- The events are mapped to WCET values
- WAF accumulates required execution time

## Arrival Monitoring and Admission

- Actual workload is calculated based on WAFs
- Check whether workload would exceed serviceable level

**Individual WAFs enforce interference bounds between real-time tasks**

Criticality Groups

**Criticality Groups**

- Tasks of one criticality level are grouped together

## Criticality Groups

- Tasks of one criticality level are grouped together
- A group of tasks is one virtual task for workload arrival

## Criticality Groups

- Tasks of one criticality level are grouped together
- A group of tasks is one virtual task for workload arrival
- Correlation of task activations improves utilization

## Criticality Groups

- Tasks of one criticality level are grouped together
- A group of tasks is one virtual task for workload arrival
- Correlation of task activations improves utilization

## Combining Monitors

## Criticality Groups

- Tasks of one criticality level are grouped together
- A group of tasks is one virtual task for workload arrival
- Correlation of task activations improves utilization

## Combining Monitors

- Monitors for groups separately

## Criticality Groups

- Tasks of one criticality level are grouped together
- A group of tasks is one virtual task for workload arrival
- Correlation of task activations improves utilization

## Combining Monitors

- Monitors for groups separately
- Each monitor respects WAF upper bounds

## Criticality Groups

- Tasks of one criticality level are grouped together
- A group of tasks is one virtual task for workload arrival
- Correlation of task activations improves utilization

## Combining Monitors

- Monitors for groups separately
- Each monitor respects WAF upper bounds
- Monitors are synchronized by guarantee interface tuples

## Criticality Groups

- Tasks of one criticality level are grouped together
- A group of tasks is one virtual task for workload arrival
- Correlation of task activations improves utilization

## Combining Monitors

- Monitors for groups separately
- Each monitor respects WAF upper bounds
- Monitors are synchronized by guarantee interface tuples
- Pareto interface tuples with maximized individual WAFs

## Properties

**Properties**

- **Event-triggered**

**Properties**

- **Event-triggered**
- **Considering worst-case is highly pessimistic**

Properties

- **Event-triggered**
- **Considering worst-case is highly pessimistic**
- **Group monitoring provides more accurate assumptions**

# Resource Management for Mixed-Criticality Systems

www.ict.tuwien.ac.at

## Resource Utilization

- Enable the derivation of tight latency bounds for guaranteed latency traffic
- Analyzable in terms of schedulability

## Resource Utilization

- Enable the derivation of tight latency bounds for guaranteed latency traffic
- Analyzable in terms of schedulability

## Quality of Service

- Reduce the average latency of best-effort traffic
- Not applicable to guaranteed latency traffic

## Resource Utilization

- Enable the derivation of tight latency bounds for guaranteed latency traffic
- Analyzable in terms of schedulability

## Quality of Service

- Reduce the average latency of best-effort traffic
- Not applicable to guaranteed latency traffic

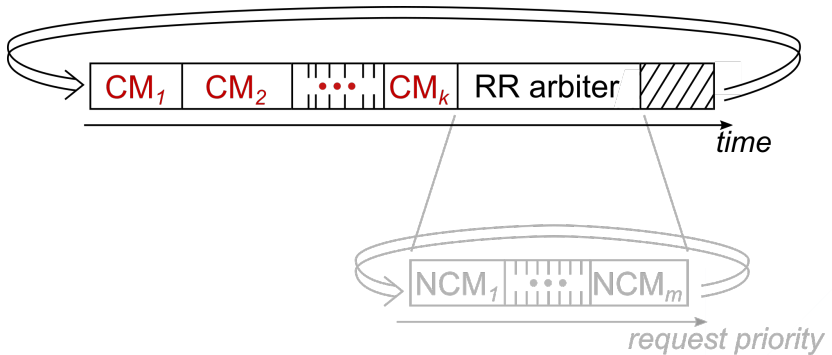## Scalability

- Scale well with the number of bus masters or NoC nodes

- First layer : TDMA for critical bus masters
- Second layer : RR for non-critical bus masters

- First layer : TDMA for critical bus masters
- Second layer : RR for non-critical bus masters

- First layer : arbitration between criticality classes
- Second layer : arbitration between tasks of the elected criticality class

- First layer : arbitration between criticality classes
- Second layer : arbitration between tasks of the elected criticality class
- Both layers are Weighted Harmonic RR

### Weighted Harmonic Round Robin

- Bus control granted for 1 request only
- Each master can get several slots per period
- Slots are evenly distributed

# CArb – Example

| $C_1$ | $C_2$ | $C_1$ | $C_3$ | $C_1$ | $C_2$ | $C_1$ | $C_3$ |
|---|---|---|---|---|---|---|---|

inter-class arbiter

| $T_{11}$ | $T_{12}$ | $T_{11}$ | $T_{13}$ |
|---|---|---|---|

class 1 arbiter

| $T_{21}$ | $T_{22}$ |
|---|---|

class 2 arbiter

| $T_{31}$ | $T_{32}$ |
|---|---|

class 3 arbiter

| $T_{11}$ | $T_{12}$ | $T_{21}$ | $T_{22}$ | $T_{21}$ | $T_{11}$ | $T_{13}$ | $T_{31}$ | $T_{11}$ | $T_{12}$ | $T_{22}$ | $T_{21}$ | $T_{22}$ | $T_{11}$ | $T_{13}$ | $T_{31}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

final schedule

$\text{WCET}_{total}(\text{T}_{13}) < exec\_time(\text{T}_{13})$
$\rightarrow$ Drop low-critical tasks, system-wide mode switch

$\text{WCET}_{total}(\text{T}_{13}) < exec\_time(\text{T}_{13})$
$\rightarrow$ Drop low-critical tasks, system-wide mode switch

$\text{WCET}_{iso}(\text{T}_{13}) < exec\_time(\text{T}_{13}) < \text{WCET}_{total}(\text{T}_{13})$
$\rightarrow$ Cut interference, arbiter-level mode switch

| $\text{T}_{11}$ | $\text{T}_{12}$ | $\text{T}_{11}$ | $\text{T}_{13}$ |
|---|---|---|---|

$\text{WCET}_{total}(\text{T}_{13}) < exec\_time(\text{T}_{13})$
$\rightarrow$ Drop low-critical tasks, system-wide mode switch

$\text{WCET}_{iso}(\text{T}_{13}) < exec\_time(\text{T}_{13}) < \text{WCET}_{total}(\text{T}_{13})$
$\rightarrow$ Cut interference, arbiter-level mode switch

| $\text{T}_{11}$ | $\text{T}_{12}$ | $\text{T}_{11}$ | $\text{T}_{13}$ |
|---|---|---|---|

If $exec\_time(\text{T}_{13}) - \text{WCET}_{iso}(\text{T}_{13}) < threshold$
$\rightarrow$ Cut only part of the interference

| $\text{T}_{11}$ | $\text{T}_{12}$ | $\text{T}_{11}$ | $\text{T}_{13}$ | $\text{T}_{21}$ | $\text{T}_{22}$ | $\text{T}_{21}$ | $\text{T}_{11}$ | $\text{T}_{12}$ | $\text{T}_{11}$ | $\text{T}_{13}$ | $\text{T}_{22}$ | $\text{T}_{21}$ | $\text{T}_{22}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

www.ict.tuwien.ac.at

www.ict.tuwien.ac.at

The header flit of critical messages contains a *Blocking Counter* field.

<div>

### Network Interface

*initialize BC value*

*(with maximum number of times the message can be delayed)*

</div>

<div>

### Router
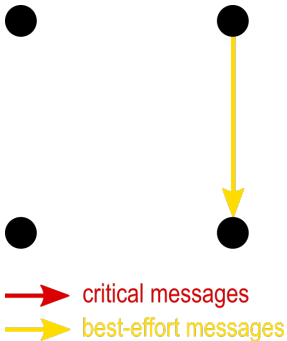
*if message is stalled*
*    decrement BC*
*    if BC = 0*
*        prioritize message*
*    end if*
*end if*

</div>

- Message BE1
  Path list : {(1,2), (2,2)}, {(1,2), (1,1), (2,1), (2,2)}

critical messages
best-effort messages

critical messages
best-effort messages

- Message BE1
  Path list : {(1,2), (2,2)}, {(1,2), (1,1), (2,1), (2,2)}

- Message C1
  Path list: {(1,2), (2,2)}

# Resource Management for Mixed-Criticality Systems

Hardware model:

- 16 Clusters

Hardware model:

- 16 Clusters
- Each containing 16 cores, 3 NoC interfaces, and 16 SRAM banks

Hardware model:

- 16 Clusters
- Each containing 16 cores, 3 NoC interfaces, and 16 SRAM banks
- Each SRAM bank has its own dedicated controller, which arbitrates between all cores and interfaces

Hardware model:

- 16 Clusters
- Each containing 16 cores, 3 NoC interfaces, and 16 SRAM banks
- Each SRAM bank has its own dedicated controller, which arbitrates between all cores and interfaces
  - NoC receiver interface always has priority
  - RR arbitration is performed between other components

Software model:

- Independent task set and scheduler for each cluster

Software model:

- Independent task set and scheduler for each cluster
- Criticality mode switch

Software model:

- Independent task set and scheduler for each cluster
- Criticality mode switch
- Scheduling data:
  - Execution times **without memory accesses**
  - Number of memory accesses
  - Dependency graph

Software model:

- Independent task set and scheduler for each cluster
- Criticality mode switch
- Scheduling data:
  - Execution times **without memory accesses**
  - Number of memory accesses
  - Dependency graph
- Tasks of different criticality cannot be scheduled on the same time frame

Software model:

- Independent task set and scheduler for each cluster
- Criticality mode switch
- Scheduling data:
    - Execution times **without memory accesses**
    - Number of memory accesses
    - Dependency graph
- Tasks of different criticality cannot be scheduled on the same time frame
- Joint computation of task schedule and data mapping to SRAM banks

When a task requires data to be fetched through the NoC

*Task A → send request for data*

*dependency constraint:*
*delay t*

*Task B → use data*

$t = 2 * \mathsf{NoC}$ worst-case latency $+$ resource worst-case response time

- A resource server "encapsulates" the access to a resource

- A resource server "encapsulates" the access to a resource
- This can allow to consider only the server and the requesting component in modeling

- A resource server "encapsulates" the access to a resource
- This can allow to consider only the server and the requesting component in modeling
- The resource itself and the interfering components can be laid aside **provided the inter-process communication protocol is adequate**
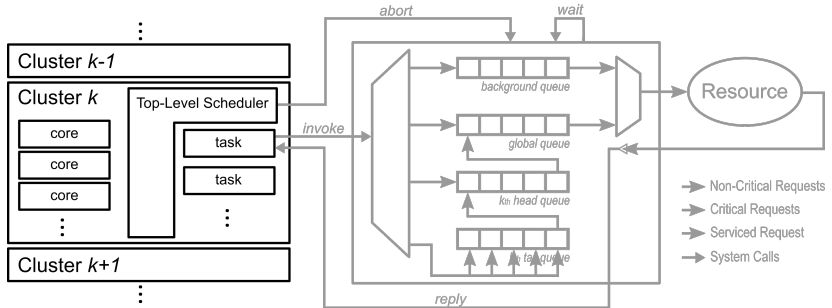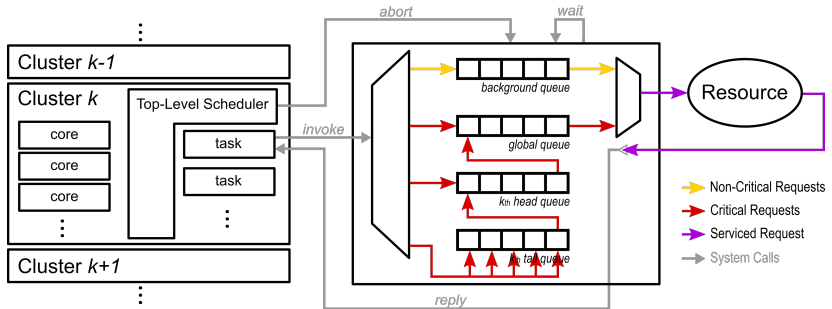
- A resource server "encapsulates" the access to a resource
- This can allow to consider only the server and the requesting component in modeling
- The resource itself and the interfering components can be laid aside **provided the inter-process communication protocol is adequate**
- For MCS, the IPC protocol should also take criticality into account

- Response to a pressing industrial need:

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?
- Mixed Criticality scheduling Theory;

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?
- Mixed Criticality scheduling Theory;
  - Design time verification;

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?
- Mixed Criticality scheduling Theory;
  - Design time verification;
  - Run-time robustness, when the system violates assumptions;

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?
- Mixed Criticality scheduling Theory;
  - Design time verification;
  - Run-time robustness, when the system violates assumptions;
- Providing sound mechanisms for sharing:

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?
- Mixed Criticality scheduling Theory;
  - Design time verification;
  - Run-time robustness, when the system violates assumptions;
- Providing sound mechanisms for sharing:
  - Processing elements;

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?
- Mixed Criticality scheduling Theory;
  - Design time verification;
  - Run-time robustness, when the system violates assumptions;
- Providing sound mechanisms for sharing:
  - Processing elements;
  - Buffers, caches, memory;

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?

- Mixed Criticality scheduling Theory;
  - Design time verification;
  - Run-time robustness, when the system violates assumptions;

- Providing sound mechanisms for sharing:
  - Processing elements;
  - Buffers, caches, memory;
  - I/O: Pins, drivers;

www.ict.tuwien.ac.at

- Response to a pressing industrial need:
  - How to provide safety and real-time guarantees in efficient implementations?

- Mixed Criticality scheduling Theory;
  - Design time verification;
  - Run-time robustness, when the system violates assumptions;

- Providing sound mechanisms for sharing:
  - Processing elements;
  - Buffers, caches, memory;
  - I/O: Pins, drivers;
  - **Communication.**

¿ Questions ?