

## 7.a WCET analysis techniques

Credits to Enrico Mezzetti, PhD  
(enrico.mezzetti@bsc.es)

### Computing the WCET /1

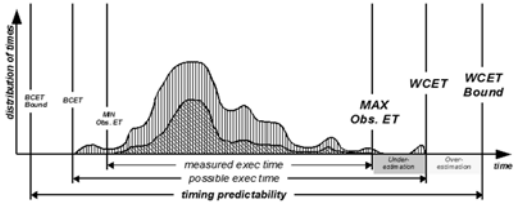
- Why not measure the WCET of a task on its real HW?
  - Worst-case input → Task
  - Worst-case HW state → Target Hardware (black box) → Logic analyser, oscilloscope, etc. **WCET?**
- Triggering the WCET by test is very difficult
  - ❑ Supplying input data that cover all possible executions of the program is an intractable problem in practice
  - ❑ Worst-case initial state is difficult to determine with modern HW
    - Complex pipelines (out-of-order execution)
    - Caches
    - Branch predictors and speculative execution

### Worst-case execution time (WCET)

- For any input data and all initial logical states
  - ❑ So that all *execution paths* of the program are covered
- For any hardware state
  - ❑ So that the worst-case *execution conditions* are in effect
- Measurement-based WCET analysis
  - ❑ On either the real HW or a cycle-accurate simulator
  - ❑ The *high-watermark* value can be  $\ll$  WCET
- Static WCET analysis
  - ❑ For an abstract model of the HW and of the program

### Computing the WCET /2

- Exact WCET not generally computable ( $\sim$  the *halting problem*)
- Yet, having WCET *bounds* is crucial to feasibility analysis
  - ❑ Must be *safe* to be an upper bound to all possible executions
  - ❑ Must be *tight* to avoid costly over-dimensioning

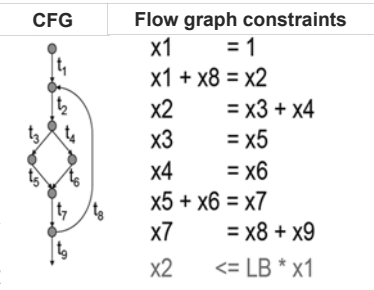


Static WCET analysis /1

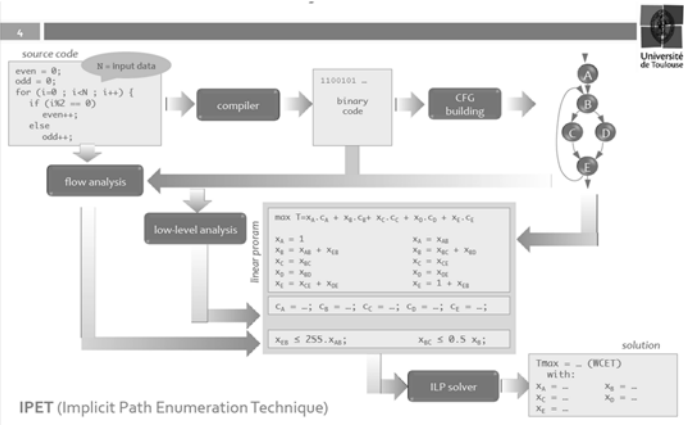
- To analyze a program without executing it
  - Needs an *abstract model* of the target HW
  - As well as the binary executable of the program
- Execution time depends on control flow and HW
  - **High-level analysis** addresses the program behavior
    - Control flow analysis builds a control flow graph (CFG)
  - **Low-level analysis** determines the timing cost of individual instructions on the abstract model of the HW
    - Not constant for modern HW
    - Must be aware of the HW inner workings (pipeline, caches, etc.)

Implicit path enumeration technique

- The program's CFG is augmented with flow graph constraints
- The WCET is computed with integer linear programming or constraint programming
- $WCET = \sum_i x_i \times t_i$ 
  - $x_i$  is the *execution frequency* of CFG edge  $i$
  - $t_i$  the *execution time* of CFG edge  $i$

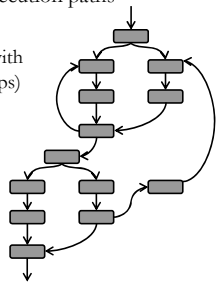


Static WCET analysis /2




Static WCET analysis /3

- **High-level analysis** /1
  - Must analyze all possible execution paths of the program
    - Builds the CFG as a superset of all possible execution paths
    - *Basic block* is the unit of that analysis
      - The longest sequence of program instructions with single entry and single exit (no branches, no loops)
  - Challenges with path analysis
    - *Input-data dependency*
    - *Infeasible paths*
    - *Loop bounds* (and recursion depth)
    - *Dynamic calls* (through pointers)



### Static WCET analysis /4

■ **High-level analysis /2**

- Several techniques are employed to allow using IPET
  - *Control-flow analysis* to construct the CFG
    - First finding the basic blocks and then building the graph among them
  - *Data-flow analysis* to find loop bounds
  - *Value analysis* to resolve memory accesses
- Automatic information extraction is insufficient 
  - User annotation of *flow facts* is needed
    - To facilitate detection of infeasible paths
    - To refine loop bounds
    - To define frequency relations between basic blocks
    - To specify the target of dynamic calls and referenced memory addresses

### Static WCET analysis /6

■ **Low-level analysis /2**

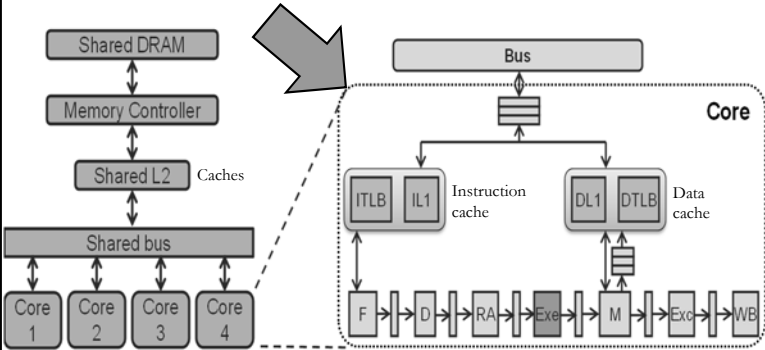
- Concrete HW states
  - Determined by the history of execution
  - Cannot compute all HW states for all possible executions
    - Invariant HW states are grouped into execution contexts
    - *Conservative overestimations* are made to reduce the research space
- *Abstract interpretation*
  - Computes abstract states and specific operators in the abstract domain
    - *Update function* to keep the abstract state current along the exec path
    - *Join function* to merge control flows after a branch
- Some techniques are specific to each HW feature

### Static WCET analysis /5

■ **Low-level analysis /1**

- Requires abstract modeling of all HW features
  - Processor, memory subsystem, buses, peripherals, ...
  - It is *conservative* : it must never underestimate actual costs
  - All possible HW states should be accounted for
- Challenges with HW modeling
  - *Precise modeling* of complex hardware is difficult
    - Inherent complexity (e.g., out-of-order pipelines)
    - Lack of comprehensive information (intellectual property, patents, ...)
    - Differences between specification and implementation (!)
  - *Exhaustive representation* of all HW states is computationally infeasible

### Understanding the hardware /1



Courtesy of **PROXIMA**

Understanding the hardware /2

The diagram illustrates the internal hardware of a processor core. It shows two main memory paths: Instruction memory and Data memory. The Instruction path includes an Instruction Store (I\$) and Instruction TLB (ITLB), with a buffer (F) before the I\$. The Data path includes a Data Store (D\$) and Data TLB (DTLB), with buffers (D, E, WB) before the D\$. The core is labeled 'Core' and the diagram is credited to 'PROXIMA'.

2016/17 UniPD / T. Vardanega

Real-Time Systems

323 of 448

Cache Associativity

The diagram explains four types of cache associativity using a bookshelf analogy. Direct Mapped: Each book (cache block) has a fixed spot. 2-Way Set Associative: Books are grouped into sets of two, and a tag helps find the set. 4-Way Set Associative: Books are grouped into sets of four, with fewer sets and index bits. Fully Associative: No index is needed; every tag must be compared. A cartoon character explains that a direct-mapped cache is just a 1-way set associative cache, and a fully associative cache is an m-way set associative cache.

2016/17 UniPD / T. Vardanega

Real-Time Systems

325 of 448

Understanding the cache

The diagram compares two mapping schemes. Direct mapping (by index) shows a 32-bit address split into a 31-bit tag and a 1-bit index (set). The cache entry has a tag, data, and a valid (V) bit. Set-associative mapping (by set) shows a 32-bit address split into a 31-bit tag and a 1-bit set. The cache entry has a tag, data, and a valid (V) bit. Both schemes show a '1. closure' and '2. look-up' process.

2016/17 UniPD / T. Vardanega

Real-Time Systems

324 of 448

Static WCET analysis: the big picture

The flowchart shows the process of Static WCET analysis. A program (exec, disassembly,...) and user annotations are input into an analysis framework and abstract hardware model. The output is safe WCET bounds.

■ Open problems

❑ Can we always trust the abstract model of the HW?

❑ How much overestimation do we incur?

- Inclusion of infeasible paths
- Overestimation is inevitable in abstract state computation

❑ Intrinsic weakness of user annotations

- Labor intensive and error prone

2016/17 UniPD / T. Vardanega

Real-Time Systems

326 of 448

## Static WCET analysis /7

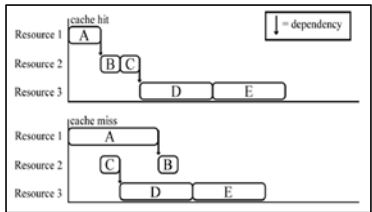
- Safeness is at risk
  - When *local* worst case does not always lead to *global* worst case
  - Which is the case when **timing anomalies** occur
    - Complex hardware architectures (e.g., out-of-order pipelines)
    - Even improper design choices (e.g., inept cache replacement policies)
    - *Counter-intuitive* timing behavior
    - Faster execution of a single instruction causes *long-term* negative effects
  - Both are very difficult to account for in static analysis

## Hybrid analysis /1

- To obtain *realistic* (less pessimistic) WCET estimates
  - On the real target processor
  - On the final executable
  - Knowing that safeness not guaranteed (!)
- Hybrid approaches exploit
  - The measurement of *basic blocks* on the real HW
    - To avoid pessimism from abstract modeling
  - Static analysis techniques to combine the obtained measures
    - Knowledge of the program execution paths
- Newer approaches explore probabilistic properties (!)

## Timing anomaly: example

- Assume some dependence between instructions
- Shared resources (e.g. pipeline stages) and opportunistic scheduling of request servicing



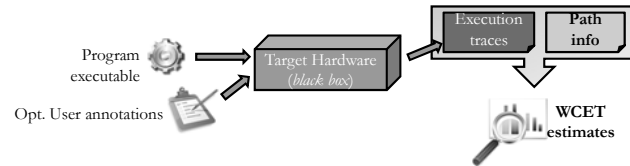
- Faster execution of A leads to a worse case overall execution owing to the order in which the instructions are executed

## Hybrid analysis /2

- Approaches to collect timing information
  - *Software instrumentation*
    - The program is augmented with instrumentation code
    - Instrumentation effects the timing behavior of the program (aka the *probe effect*) and causes problems to deciding what's the final system
  - *Hardware instrumentation*
    - Depends on specialized HW features (e.g., debug interface)
- Confidence in the results contingent on the coverage of the executions and on the exploration of worst-case states
  - Exposed to the same problems as static analysis and measurement
  - *Worst-case state dependence is gone if HW response time is randomized*



## Hybrid analysis: the big picture



- Open problems
  - Can we trust the resulting estimates?
    - Contingent on worst-case input and worst-case HW state
    - Consideration of infeasible paths
  - Needs the real execution environment or an identical copy of it
    - May cause serious cost impact and inherent difficulty of exactness

## Selected readings

- R. Wilhelm et al. (2008)  
*The worst-case execution-time problem—overview of methods and survey of tools*  
 DOI: 10.1145/1347375.1347389

## Summary

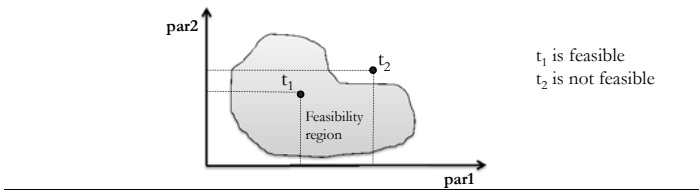
- The challenge of computing the WCET
- Static analysis
  - High-level analysis
  - Low-level analysis
- Hybrid analysis (measurement-based)

## 7.b Schedulability analysis techniques

Credits to Marco Panunzio, PhD  
 (marco.panunzio@thalesaleniaspace.com)

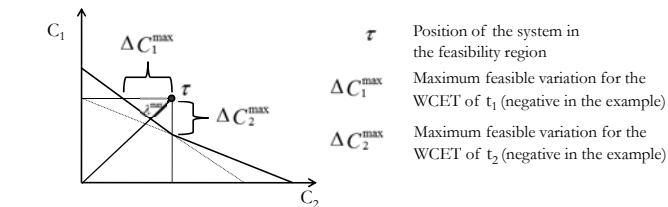
### Feasibility region

- The topological space that represents the set of feasible systems with respect to the workload model parameters
  - N-dimensional space with N-parameter analysis
  - Function of the timing parameters
  - Specific to the scheduling algorithm in use



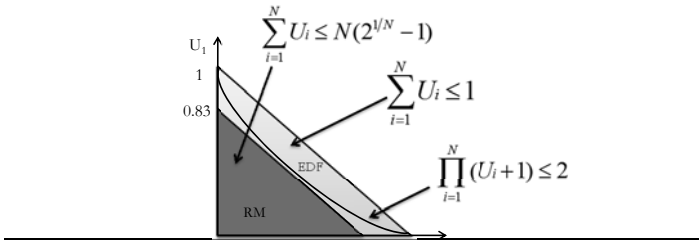
### Sensitivity analysis /1

- Investigates the changes in a given system that
  - Improve the fit of an already feasible system
  - Make feasible an infeasible system



### Advanced utilization tests

- *Hyperbolic bound* improves Liu & Layland utilization test
  - For systems with periodic tasks under FPS and DMPO
  - E. Bini and G. Buttazzo: “*A Hyperbolic Bound for the Rate Monotonic Algorithm*”. Proceedings of the 13<sup>th</sup> ECRTS, 2001



### Sensitivity analysis /2

- Major computational complexity
- Theory still under development
  - Does not account for shared resources, multi-node systems, partitioned systems
- Interesting potential
  - To explore solution space in the *dimensioning* phase of design
    - Presently only applicable to period/MIAT and WCET
  - To study the consequences of changes to timing parameters
    - To allow for the inclusion of better functional value in the system
    - To renegotiate timing (or functional) parameters

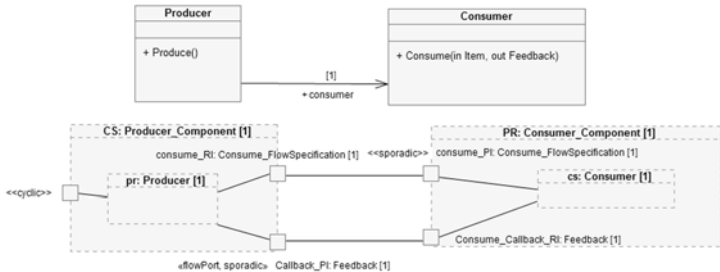
Transactions /1

- Causal relations between activities
  - They consider information relevant to analysis that is not captured by classic workload models
    - Dependence in the activation of jobs
- Originally introduced for the analysis of distributed systems
  - Also useful for the analysis of “collaboration patterns” employed for single-CPU systems



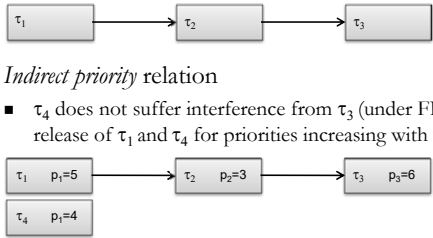
Example /1

- A “callback pattern” to permit **in out** interactions between tasks in Ravenscar systems

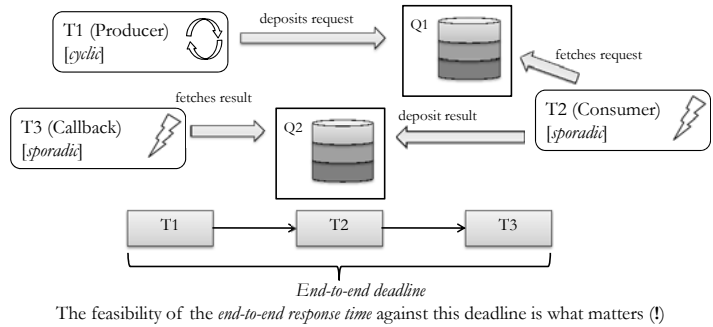


Transactions /2

- Two main kinds of dependence
  - Direct precedence relation (e.g., producer-consumer)
    - $\tau_2$  cannot proceed until  $\tau_1$  completes
- Indirect priority relation
  - $\tau_4$  does not suffer interference from  $\tau_3$  (under FPS and synchronous release of  $\tau_1$  and  $\tau_4$  for priorities increasing with values)



Example /2





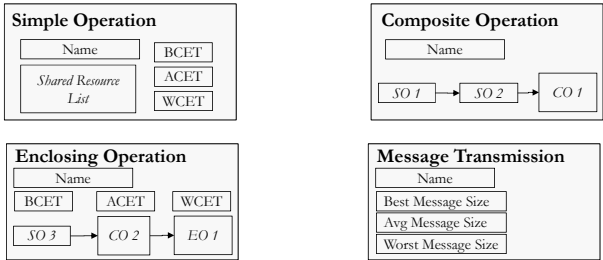
## 343 of 448



## 344 of 448

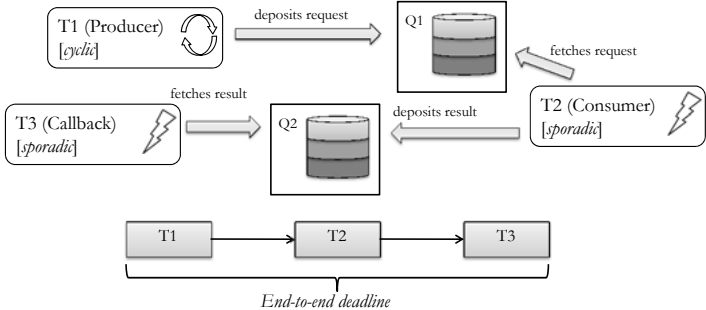
## 346 of 448

MAST: operations

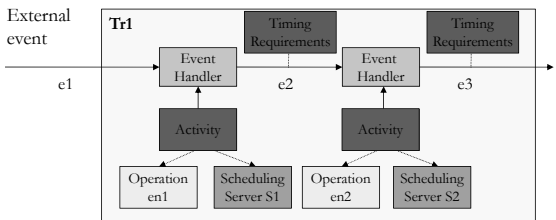


- The real-time model includes the description of all the operations in the system

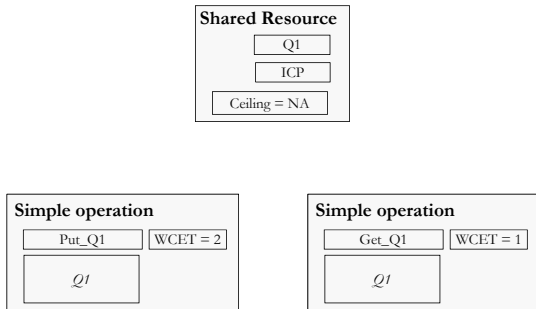
Example: Ravenscar callback



MAST: creation of a transaction



Example: shared resources in MAST



Example: modeling tasks in MAST

Simple operation

Produce\_SO

WCET = 8

None

Enclosing Operation

Produce\_EO

WCET=10

Produce\_SO

Put\_Q1

Scheduling Server

Producer\_SS

FPP Priority = 4

CPU1.P5

External event

E1

T=40

Producer

Event Handler

D = 40

Activity

Operation Produce\_EO

Scheduling Server Producer\_SS

2016/17 UniPD / T. VardanegaReal-Time Systems351 of 448

Example: classic RTA results

Producer [1] (C) T<sub>1</sub>=40 C<sub>1</sub>=10 p<sub>1</sub>=4

Consumer [2] (S) T<sub>2</sub>=40 C<sub>2</sub>=10 p<sub>2</sub>=2

Callback [3] (S) T<sub>3</sub>=40 C<sub>3</sub>=5 p<sub>3</sub>=5

Q1 Ceiling=4

Q2 Ceiling=5

B<sub>1</sub>=2 B<sub>2</sub>=0 B<sub>3</sub>=2

Classic RTA

R<sub>1</sub> = 17

R<sub>2</sub> = 25

R<sub>3</sub> = 7

This misses out completely that T<sub>3</sub> is to be preceded by T<sub>2</sub> and T<sub>1</sub> (!)

2016/17 UniPD / T. VardanegaReal-Time Systems353 of 448

Example: timing attributes

Producer [1] (C) T<sub>1</sub>=40 C<sub>1</sub>=10 p<sub>1</sub>=4

Consumer [2] (S) T<sub>2</sub>=40 C<sub>2</sub>=10 p<sub>2</sub>=2

Callback [3] (S) T<sub>3</sub>=40 C<sub>3</sub>=5 p<sub>3</sub>=5

Q1 Ceiling=4

Q2 Ceiling=5

2016/17 UniPD / T. VardanegaReal-Time Systems352 of 448

Example: introducing transactions

External event

E1

T=40

Producer\_TR

Event Handler

D = 40

Activity

Operation Produce\_EO

Scheduling Server Producer\_SS

Event Handler

D = 40

Activity

Operation Consume\_EO

Scheduling Server Consumer\_SS

Event Handler

D = 40

Activity

Operation Callback\_EO

Scheduling Server Callback\_SS

2016/17 UniPD / T. VardanegaReal-Time Systems354 of 448

# Example: end-to-end analysis

<b>Producer</b> [1]	(C)	$T_1=40$	$C_1=10$	$p_1=4$
<b>Consumer</b> [2]	(S)	$T_2=40$	$C_2=10$	$p_2=2$
<b>Callback</b> [3]	(S)	$T_3=40$	$C_3=5$	$p_3=5$

**Q1**    Ceiling=4  
**Q2**    Ceiling=5

$\Rightarrow$      $B_1=2$      $B_2=0$      $B_3=2$

<b>Classic RTA</b>	<b><i>Precedence and offset-based</i></b>	
$R_1 = 17$	$R_1(Ir) = 12$	
$R_2 = 25$	$R_2(Ir) = 20$	$\leftarrow$ <i>Response time relative to the beginning of the transaction!</i>
$R_3 = 7$	$R_3(Ir) = 27$	

# Summary

- Feasibility region
- Advanced utilization tests
- Sensitivity analysis
- Transactions
- Example tool (MAST)