

INSTITUTE OF COMMUNICATION, INFORMATION AND PERCEPTION TECHNOLOGIES
Scuola Superiore Sant'Anna

Retis
Real-Time Systems Laboratory

Response-Time Analysis of Conditional DAG Tasks in Multiprocessor Systems

Alessandra Melani

Retis 1

What does it mean?

- ❑ « Response-time analysis » ✓
- ❑ « conditional »
- ❑ « DAG tasks »
- ❑ « multiprocessor systems » ✓

Retis 2

What does it mean?

- ❑ « Response-time analysis »
- ❑ « conditional »
- ❑ « DAG tasks »
- ❑ « multiprocessor systems »

If-then-else statements Switch statements

Retis 3

What does it mean?

- ❑ « Response-time analysis »
- ❑ « conditional »
- ❑ « DAG tasks »
- ❑ « multiprocessor systems »

DAG: Directed Acyclic Graph

Retis 4

In other words

- ❑ We will analyze a **multiprocessor** real-time systems...
- ❑ ... by means of a **schedulability test** based on **response-time analysis**
- ❑ ... assuming **Global Fixed Priority** or **Global EDF** scheduling policies
- ❑ ... and assuming a **parallel task model** (i.e., a task is modelled as a **Directed Acyclic Graph - DAG**)

Retis 5

Parallel task models

Many parallel programming models have been proposed to support parallel computation on multiprocessor platforms (e.g., OpenMP, OpenCL, Cilk, Cilk Plus, Intel TBB)

OpenMP

OpenCL

Cilk Plus

TBB

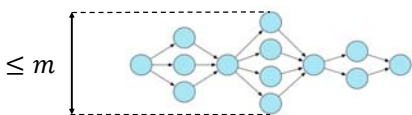
Early real-time scheduling models: each recurrent task is completely sequential

Recently, more expressive execution models allow exploiting task parallelism

Retis 6

Fork-join

- Each task is an alternating sequence of sequential and parallel segments
- Every parallel segment has a degree of parallelism $\leq m$ (number of processors)



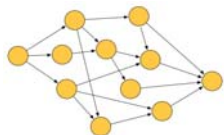
Synchronous-parallel

- Generalization of the fork-join model
- Allows consecutive parallel segments
- Allows an arbitrary degree of parallelism of every segment
- Synchronization at segment boundaries: a sub-task in the new segment may start only after completion of all sub-tasks in the previous segment



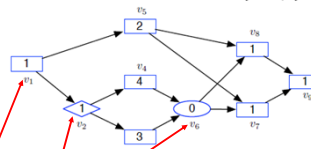
DAG

- Directed acyclic graph (DAG) $G_i = (V_i, E_i)$
- $V_i = \{v_{i,1}, \dots, v_{i,n_i}\}; E_i \subseteq V_i \times V_i$
- Generalization of the previous two models
- Every node is a sequential sub-task
- Arcs represent precedence constraints between sub-tasks



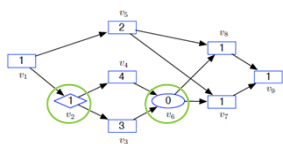
cp-DAG

- Conditional - parallel DAG (cp-DAG) $G_i = (V_i, E_i)$



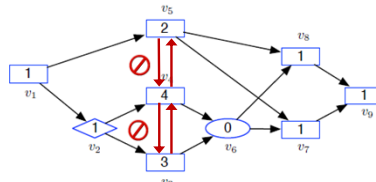
- Two types of nodes
 - Regular:** all successors must be executed in parallel
 - Conditional:** to model start/end of a conditional construct (e.g., if-then-else statement)
- Each node has a WCET $C_{i,j}$
- In this lecture, we will focus on **this** task model

Conditional pairs



- (v_2, v_6) form a **conditional pair**
 - v_2 is a starting conditional node
 - v_6 is the joining point of the conditional branches starting at v_2
- Restriction:** there cannot be any connection between a node belonging to a branch of a conditional statement (e.g., v_4) and nodes outside that branch (e.g., v_5), including other branches of the same statement

Why this restriction?

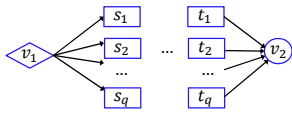


- It does not make sense for v_5 to wait for v_4 if v_3 is executed
- Analogously, v_4 cannot be connected to v_3 since only one is executed
- Violation of the correctness of conditional constructs and the semantics of the precedence relation

Formal definition (1)

Let (v_1, v_2) be a pair of conditional nodes in a DAG $G_i = (V_i, E_i)$.
The pair (v_1, v_2) is a conditional pair if the following hold:

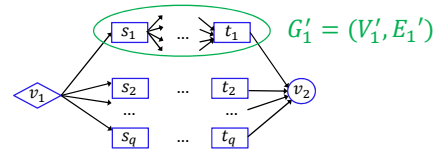
- Suppose there are exactly q outgoing arcs from v_1 to the nodes s_1, s_2, \dots, s_q , for some $q > 1$. Then there are exactly q incoming arcs into v_2 in E_i , from some nodes t_1, t_2, \dots, t_q



Formal definition (2)

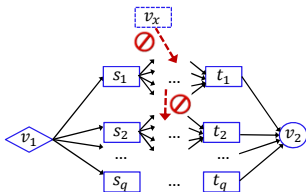
- For each $l \in \{1, 2, \dots, q\}$, let $V'_l \subseteq V_i$ and $E'_l \subseteq E_i$ denote all the nodes and arcs on paths reachable from s_l that do not include v_2 .

By definition, s_l is the sole source node of the DAG $G'_l = (V'_l, E'_l)$. It must hold that t_l is the sole sink node of G'_l .



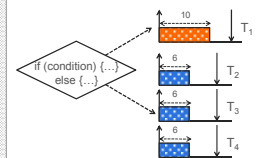
Formal definition (3)

- It must hold that $V'_l \cap V'_j = \emptyset$ for all $l, j, l \neq j$.
Additionally, with the exception of (v_1, s_l) , there should be no arcs in E_i into nodes in V'_l from nodes not in V'_l , for each $l \in \{1, 2, \dots, q\}$.
That is, $E_i \cap ((V_i \setminus V'_l) \times V'_l) = \{(v_1, s_l)\}$ should hold for all l .



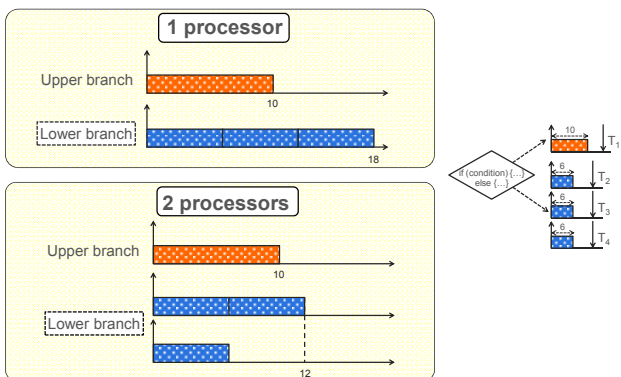
How is parallel code structured?

```
#pragma omp parallel num_threads(N)
{
  #pragma omp master {
    #pragma omp task // T0
    if (condition) {
      #pragma omp task // T1
    } else {
      #pragma omp task // T2
      #pragma omp task // T3
      #pragma omp task // T4
    }
  }
}
```

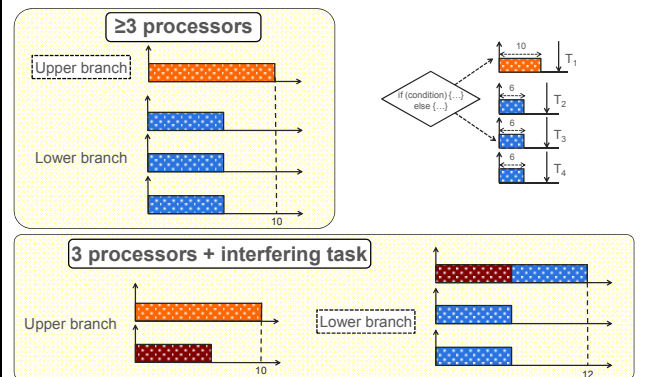


Which branch leads to the worst-case response-time?

Which branch leads to the WCRT?



Which branch leads to the WCRT?



Lesson learnt

Depending on the number of processors and on the interfering tasks, it is not obvious to identify the **branch leading to the WCRT**

It makes sense to account for the different execution flows by **enriching the task model**

Why don't we do it also with **sequential tasks**?

- Only the longest path matters
- Conditional branches are already incorporated in the notion of WCET

19

System model

- n conditional-parallel tasks (cp-tasks) τ_i , expressed as cp-DAGs in the form $G_i = (V_i, E_i)$
- platform composed of m identical processors
- **sporadic** arrival pattern (minimum inter-arrival time T_i between jobs of task τ_i)
- **constrained** relative deadline $D_i \leq T_i$

Problem

Schedulability analysis for cp-tasks, **globally** scheduled on **m identical processors with any work-conserving algorithm (including G-FP and G-EDF)**

20

Quantities of interest

1. Chain (or path) of a cp-task
2. Longest path
3. Volume
4. Worst-case workload
5. Critical chain

21

1. Chain (or path)

A chain (or path) of a cp-task τ_i is a sequence of nodes $\lambda = (v_{i,a}, \dots, v_{i,b})$ such that $(v_{i,j}, v_{i,j+1}) \in E_i, \forall j \in [a, b)$.

22

1. Chain (or path)

A chain (or path) of a cp-task τ_i is a sequence of nodes $\lambda = (v_{i,a}, \dots, v_{i,b})$ such that $(v_{i,j}, v_{i,j+1}) \in E_i, \forall j \in [a, b)$.

The length of the chain, denoted by $len(\lambda)$, is the sum of the WCETs of all its nodes:

$$len(\lambda) = \sum_{j=a}^b C_{i,j}$$

23

2. Longest path

The longest path L_i of a cp-task τ_i is any source-sink chain of the task that achieves the longest length

L_i also represents the time required to execute it when the number of processing units is infinite (large enough to allow maximum parallelism)

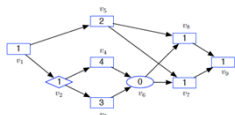
Necessary condition for feasibility: $L_i \leq D_i$

24

2. Longest path

How to compute the longest path?

- Find a topological order of the given cp-DAG
 - A topological order is such that if there is an arc from u to v in the cp-DAG, then u appears before v in the topological order → can be done in $O(n)$
 - Example: for this cp-DAG possible topological orders are
 - $(v_1, v_2, v_5, v_3, v_4, v_6, v_8, v_7, v_9)$
 - $(v_1, v_5, v_2, v_3, v_4, v_6, v_7, v_8, v_9)$
 - $(v_1, v_2, v_4, v_3, v_6, v_5, v_8, v_7, v_9)$



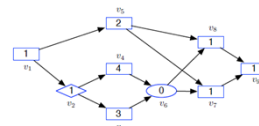
2. Longest path

How to compute the longest path?

- For each vertex $v_{i,j}$ of the cp-DAG in the topological order, compute the length of the longest path ending at $v_{i,j}$ by looking at its incoming neighbors and adding $C_{i,j}$ to the maximum length recorded for those neighbors
- If $v_{i,j}$ has no incoming neighbors, set the length of the longest path ending at $v_{i,j}$ to $C_{i,j}$

Example:

- For v_1 , record 1
- For v_2 , record 2
- For v_3 , record 5
- For v_4 , record 6
- For v_5 , record $\max(5, 6) = 6$

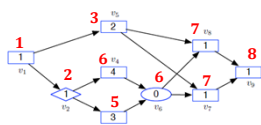


2. Longest path

How to compute the longest path?

- Finally, the longest path in the cp-DAG may be obtained by starting at the vertex $v_{i,j}$ with the largest recorded value, then repeatedly stepping backwards to its incoming neighbor with the largest recorded value, and reversing the sequence found in this way

Example: **recorded values**



- Starting at v_9 and stepping backward we find the sequence $(v_9, v_7, v_6, v_4, v_2, v_1)$
- The longest path is then $(v_1, v_2, v_4, v_6, v_7, v_9)$

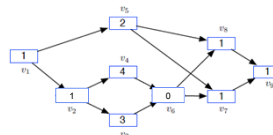
Complexity of the longest path computation: $O(n)$

3. Volume

In the **absence** of conditional branches, the volume of a task is the worst-case execution time needed to complete it on a dedicated single-core platform

It can be computed as the sum of the WCETs of all its vertices:

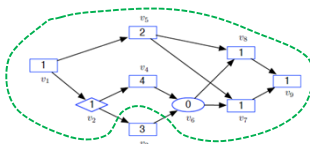
$$vol_i = \sum_{v_{i,j} \in V_i} C_{i,j}$$



It also represents the maximum amount of workload generated by a single instance of a DAG-task

4. Worst-case workload

In the **presence** of conditional branches, the worst-case workload of a task is the worst-case execution time needed to complete it on a dedicated single-core platform, over all combination of choices for the conditional branches

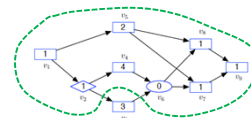


It also represents the maximum amount of workload generated by a single instance of a cp-task

In this example, the worst-case workload is given by all the vertices except v_3 , since the branch corresponding to v_4 yields a larger workload

4. Worst-case workload

How can it be computed?



Algorithm 1 Worst-Case Workload Computation

```

1: procedure WCV(G)
2:    $\sigma \leftarrow \text{TOPOLOGICAL\_ORDER}(G)$ 
3:   for  $i \leftarrow |V|$  down to 1 do
4:      $j \leftarrow \sigma(i)$  // takes the  $i^{\text{th}}$  element of the permutation
5:      $S(v_j) \leftarrow \{v_j\}$  //  $S$  takes the accumulated worst-case workload from  $v_j$  till the end of the cp-DAG
6:     if  $\text{SUCCESSORS}(v_j) \neq \emptyset$  then // if the vertex has some successors
7:       if  $\text{ISBEGINCOND}(v_j)$  then // if the vertex is the head node of a conditional pair
8:          $v^* \leftarrow \text{argmax}_{v \in \text{SUCCESSORS}(v_j)} C(S(v))$  //  $v^*$  is the successor of  $v_j$  achieving the largest partial workload
9:          $S(v_j) \leftarrow S(v_j) \cup S(v^*)$  //  $S(v^*)$  is merged into  $S(v_j)$ 
10:      else // if instead the vertex is a regular one
11:         $S(v_j) \leftarrow S(v_j) \cup \bigcup_{v \in \text{SUCCESSORS}(v_j)} S(v)$  // the workload contribution of all successors is merged into  $S(v_j)$ 
12:      end if
13:    end if
14:  end for
15:  return  $C(S(v_{\sigma(1)}))$  // the worst-case workload accumulated by the source vertex is returned as output
16: end procedure
    
```

4. Worst-case workload

What is the complexity of this algorithm?

```

Algorithm 1 Worst-Case Workload Computation
1: procedure WCW( $G$ )
2:    $\sigma \leftarrow \text{TOPOLOGICALORDER}(G)$ 
3:   for  $z = |V|$  down to 1 do
4:      $i \leftarrow \sigma(z)$ 
5:      $S(v_i) \leftarrow \{v_i\}$ 
6:     if  $\text{ISBEGINCOND}(v_i) \neq \emptyset$  then
7:       if  $\text{ISBEGINCOND}(v_i)$  then
8:          $v^* \leftarrow \text{argmax}_{v \in \text{succ}(v_i)} C(S(v))$ 
9:          $S(v_i) \leftarrow S(v_i) \cup S(v^*)$ 
10:      else
11:         $S(v_i) \leftarrow S(v_i) \cup \bigcup_{v \in \text{succ}(v_i)} S(v)$ 
12:      end if
13:    end if
14:  end for
15:  return  $C(S(v_{\sigma(1)}))$ 
16: end procedure
    
```

- $O(|E|)$ set operations
- Any of them may require to compute $C(S(v_i))$, which has cost $O(|V|)$

The time complexity is then $O(|E||V|)$

5. Critical chain

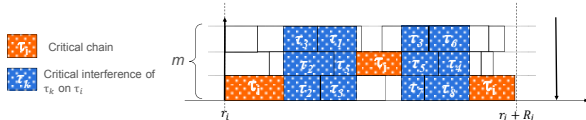
- Given a set of cp-tasks and a (work-conserving) scheduling algorithm, the **critical chain** λ_i^* of a cp-task τ_i is the chain of vertices of τ_i that leads to its worst-case response-time R_i
- How can it be identified?
 - We should know the worst-case instance of τ_i (i.e., the job of τ_i that has the largest response-time in the worst-case scenario)
 - Then we should take its sink vertex v_{i,m_i} and recursively pre-pend the last to complete among the predecessor nodes, until the source vertex $v_{i,1}$ has been included in the chain

Key observation: the critical chain is unknown, but is always upper-bounded by the longest path of the cp-task!

Critical interference

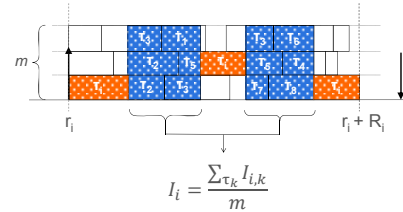
To find the response-time of a cp-task, it is sufficient to characterize the maximum interference suffered by its critical chain

The **critical interference** $I_{i,k}$ imposed by task τ_k on task τ_i is the cumulative workload executed by vertices of τ_k while a node belonging to the critical chain of τ_i is ready to execute but is not executing



Work-conserving schedulers

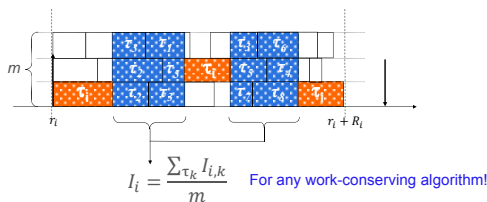
Global schedulers are typically work-conserving (e.g., Global FP/EDF)
Property: a ready job cannot execute only if all m processors are busy



We can safely assume that the interference is distributed across all m processors

Critical interference

- I_i : total interference suffered by task τ_i
- $I_{i,k}$: total interference of task τ_k on task τ_i

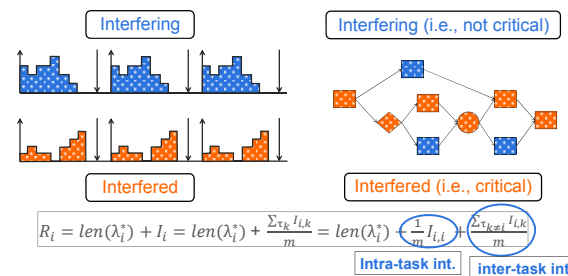


$$R_i = \text{len}(\lambda_i^*) + I_i = \text{len}(\lambda_i^*) + \frac{\sum_{\tau_k} I_{i,k}}{m}$$

Types of interference

We need to deal with two types of interference:

- Inter-task interference:** from other tasks in the system; analogous to the classic notion
- Intra-task interference:** from vertices of the same task on itself; peculiar to parallel tasks only



$$R_i = \text{len}(\lambda_i^*) + I_i = \text{len}(\lambda_i^*) + \frac{\sum_{\tau_k} I_{i,k}}{m} = \text{len}(\lambda_i^*) + \frac{1}{m} I_{i,i} + \frac{\sum_{k \neq i} I_{i,k}}{m}$$

Intra-task int. inter-task int.

Inter-task interference

- Caused by other cp-tasks executing in the system
- Finding it exactly is difficult
- We need to find an **upper-bound on the workload** of an interfering task in the scheduling window $[r_i, r_i + R_i]$
- In the sequential case (global multiprocessor scheduling):

What is the scenario that maximizes the interfering workload?

etis 37

Inter-task interference

- Sequential case**
 - The first job of τ_k starts executing as late as possible, with a starting time aligned with the beginning of the scheduling window
 - Later jobs are executed as soon as possible
- Parallel case**
 - This scenario may not give a safe upper-bound on the interfering workload. Why?

Shifting right the scheduling window may give a larger interfering workload!

etis 38

Inter-task interference

- Pessimistic assumption**
 - Each interfering job of task τ_k executes for its worst-case workload W_k
 - The carry-in and carry-out contributions are evenly distributed among all m processors
 - Distributing them on less processors cannot increase the workload within the window
 - Other task configurations cannot lead to a higher workload within the window

etis 39

Inter-task interference

- Lemma:** An upper-bound on the workload of an interfering task τ_k in a scheduling window of length L is given by

$$W_k(L) = \left\lfloor \frac{L + R_k - W_k/m}{T_k} \right\rfloor W_k + \min \left(W_k, m \cdot \left((L + R_k - \frac{W_k}{m}) \bmod T_k \right) \right)$$

- Proof:**
 - The maximum number of carry-in and body instances within the window is

etis 40

Inter-task interference

- Proof (continued):** $W_k(L) = \left\lfloor \frac{L + R_k - W_k/m}{T_k} \right\rfloor W_k + \min \left(W_k, m \cdot \left((L + R_k - \frac{W_k}{m}) \bmod T_k \right) \right)$
- Each of the $\left\lfloor \frac{L + R_k - W_k/m}{T_k} \right\rfloor$ instances contributes for W_k
- The portion of the carry-out job included in the window is $(L + R_k - \frac{W_k}{m}) \bmod T_k$

- At most m processors may be occupied by the carry-out job
- The carry-out job cannot execute for more than W_k units

etis 41

Intra-task interference

It is the interference from vertices of the same task on itself

Who is **interfering** and who is **interfered**?

- The **interfered** contribution is the **critical chain**
- Critical chain:** chain that leads to the WCRT of the cp-task

Critical chain \neq longest path

- Longest path is 10 time-units
- Critical chain can be either 10 or 6

etis 42

Thank you!

Alessandra Melani
alessandra.melani@sssup.it

etis
49