# 7.b Seeking the lost optimality

Greg Levin[†]    Shelby Funk[‡]    Caitlin Sadowski[†]

Ian Pye[†]    Scott Brandt[†]

[†]University of California
Santa Cruz

[‡]University of Georgia
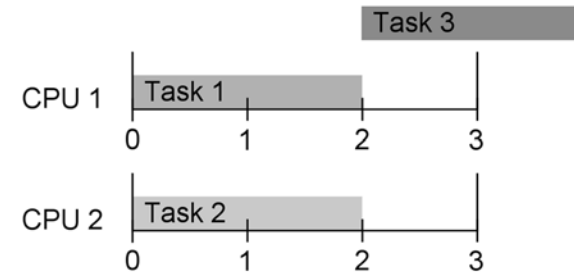Athens

---

## Partitioned Schedulers ≠ Optimal

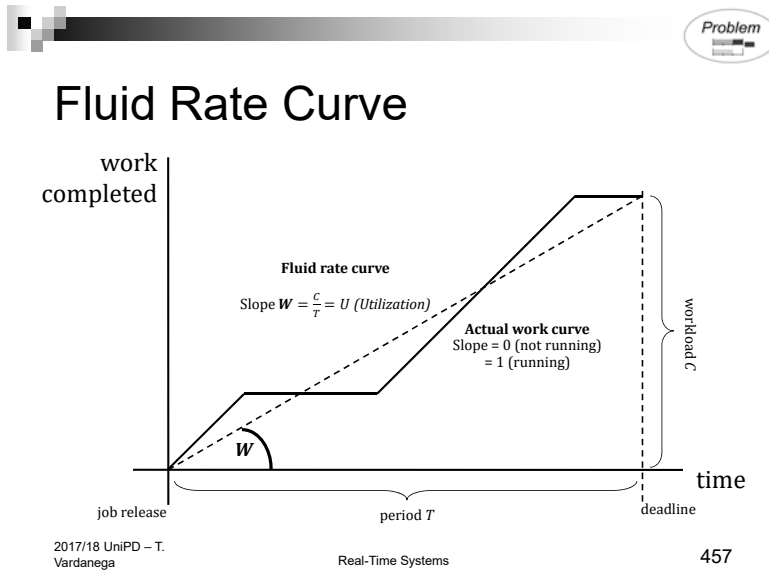■ Example: 2 processors; 3 tasks, each with 2 units of work required every 3 time units: (3,2)

---

## Global Schedulers May Succeed

■ Example: 2 processors; 3 tasks, each with 2 units of work required every 3 time units

Task 3 *migrates* between processors

## Fluid Rate Curve



Fluid rate curve

Slope $W = \frac{C}{T} = U$ (Utilization)

**Actual work curve**
Slope = 0 (not running)
= 1 (running)

workload $C$

work completed

$W$

job release        period $T$        deadline        time

## Feasible Work Region



work complete

CPU rate = 1

zero laxity

workload $C$

work completed

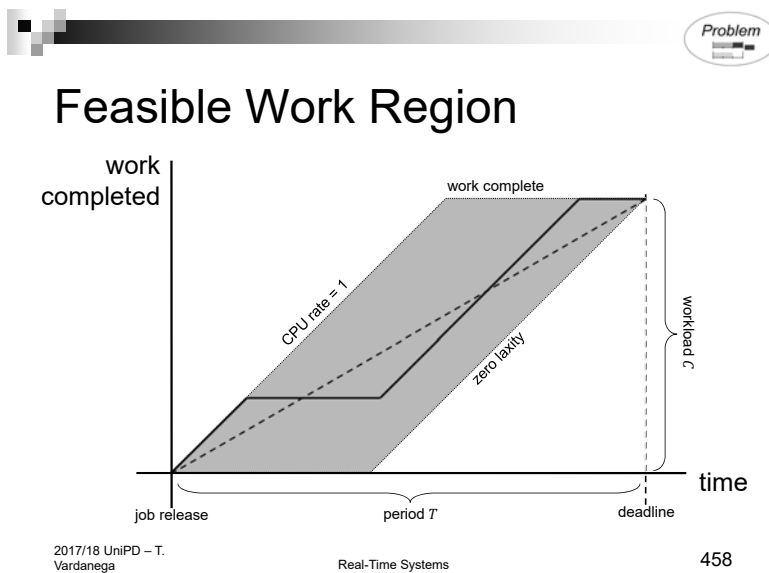job release        period $T$        deadline        time

## The Grand Challenge (*Mark 1*)

- Design an *optimal* scheduling algorithm for periodic task sets on *multiprocessors*

  □ A task set is *feasible* if there exists a schedule that meets all deadlines

  □ A scheduler is *optimal* if it can always schedule any feasible task set

## Necessary and Sufficient Conditions

- Any set of (independent) tasks needing at most
  □ 1 processor for each task $\tau_i$ ($\forall i\ U_i \leq 1$)
  □ $m$ processors for all tasks ($\sum_i U_i \leq m$)
  is feasible
- **Proof**: small scheduling intervals can approximate the fluid rate curve (at what cost?)
  □ **Status**: solved. P-Fair (1996) was the first optimal algorithm

Real Time Systems          2

## The Grand Challenge (*Mark 2*)

- Design an *optimal* scheduling algorithm with *fewer* context switches and migrations

  - Finding a feasible schedule with *the fewest* migrations is NP-Complete!

---

## Why Greedy Algorithms Fail On Multiprocessors

*DP-Fair*

- Example $(n = 3, m = 2)$

**Task 1** : Work = 9 , Period = 10
**Task 2** : Work = 9 , Period = 10
**Task 3** : Work = 8 , Period = 40

0          10                    40

**Utilization**: 9/10 + 9/10 + 8/40 = 2

---

## The Grand Challenge (*Mark 2*)

*Problem*

- Design an *optimal* scheduling algorithm with *fewer* context switches and migrations

- Status: *Solved*

  - **BUT** the solutions are complex and confusing

- **Our Contributions:** A *simple, unifying theory* for optimal global multiprocessor scheduling *and* a simple optimal algorithm
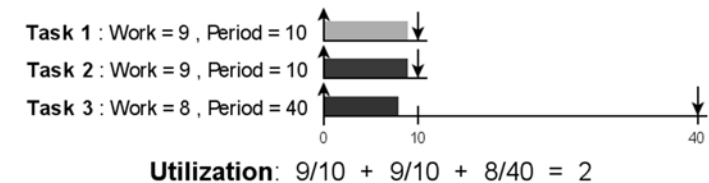
---

## Why Greedy Algorithms Fail On Multiprocessors

*DP-Fair*

At $t = 0$, $\tau_1, \tau_2$ are the obvious greedy choice

**Task 1** : Work = 9 , Period = 10
**Task 2** : Work = 9 , Period = 10
**Task 3** : Work = 8 , Period = 40

0          10                    40

CPU 1

CPU 2

0  1  2  3  4  5  6  7  8  9  10  11  12

DP-Fair
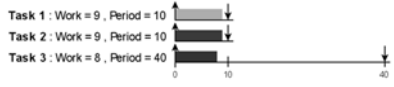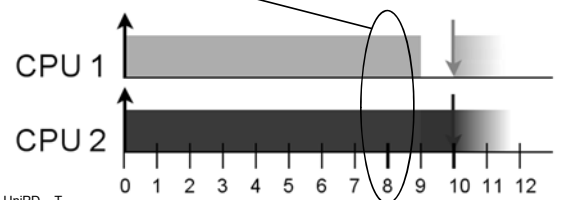
## Why Greedy Algorithms Fail On Multiprocessors

Even at $t = 8$, $\tau_1, \tau_2$ are the only "reasonable" greedy choice

Task 1 : Work = 9 , Period = 10
Task 2 : Work = 9 , Period = 10
Task 3 : Work = 8 , Period = 40

---

DP-Fair

## Why Greedy Algorithms Fail On Multiprocessors

How can we "see" this critical event at $t = 8$?

Task 1 : Work = 9 , Period = 10
Task 2 : Work = 9 , Period = 10
Task 3 : Work = 8 , Period = 40

---

DP-Fair

## Why Greedy Algorithms Fail On Multiprocessors

Yet, if $\tau_3$ isn't started by $t = 8$, the resultant idle time eventually causes a deadline miss

Task 1 : Work = 9 , Period = 10
Task 2 : Work = 9 , Period = 10
Task 3 : Work = 8 , Period = 40

---

DP-Fair

## **Proportioned** Algorithms Succeed On Multiprocessors

Subdivide $\tau_3$ in two subtasks with the same period as $\tau_1, \tau_2$

Task 1 : Work = 9 , Period = 10
Task 2 : Work = 9 , Period = 10
Task 3 ( Work = 2 , Period = 10 )

## Proportioned Algorithms Succeed On Multiprocessors

Now $\tau_3$ has a *zero-laxity* event at $t = 8$

Task 1 : Work = 9 , Period = 10
Task 2 : Work = 9 , Period = 10
Task 3 : Work = 2 , Period = 10

CPU 1

CPU 2

0 1 2 3 4 5 6 7 8 9 10 11 12

2017/18 UniPD – T. Vardanega

Real-Time Systems

469

---

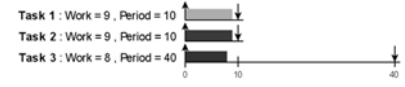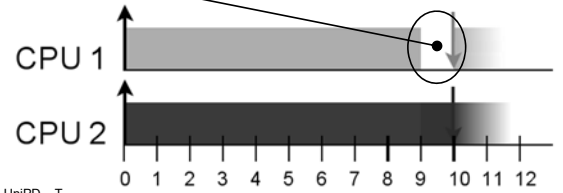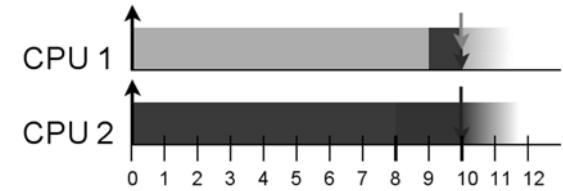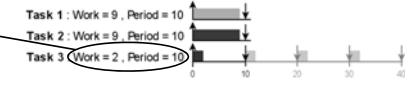## Proportional Fairness

- **Insight:** scheduling is easier when all jobs have the same deadline

> Theorem [Hong, Leung: RTSS 1988, IEEE TCO 1992]
> *No optimal on-line scheduler can exist for a set of jobs with two or more distinct deadlines on any m multiprocessor system, where m > 1*
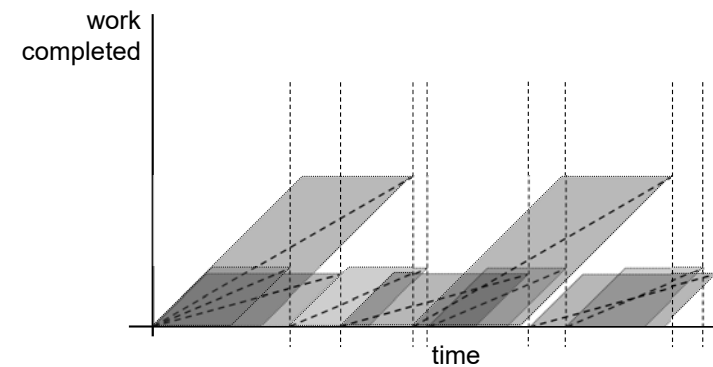
- **Application:** apply all deadlines to all jobs
  - Assign workloads proportional to utilization
  - Work complete matches fluid rate curve at every system deadline

2017/18 UniPD – T. Vardanega

Real-Time Systems

470

---

## Proportional Fairness is the Key

- All known optimal algorithms enforce proportional fairness at all deadlines
  - **P-Fair** (1996) - *Baruah, Cohen, Plaxton, and Varvel*
    ( the extreme: proportional fairness at *all times* )
  - **BF** (2003) - *Zhu, Mossé, and Melhem*
  - **LLREF** (2006) - *Cho, Ravindran, Jensen*
  - **EKG** (2006) - *Andersson, Tovar*

- Why do they all use proportional fairness?

2017/18 UniPD – T. Vardanega

Real-Time Systems

471

---

## Scheduling Multiple Tasks is Complicated

work completed

time

2017/18 UniPD – T. Vardanega

Real-Time Systems

472

## Scheduling Multiple Tasks with Same Deadline is Easy

work completed

time

## Restricted Feasible Regions Under Deadline Partitioning

work completed

all system deadlines

time

## Actual Feasible Regions

work completed

job deadlines

time

## The DP-Fair Scheduling Policy

- Partition time into *slices* based on all system deadlines

- Allocate each job a per-slice workload equal to its *utilization* times the *length of the slice*

- Schedule jobs within each slice in any way that obeys the following three rules:

  1. Always run a job with zero *local laxity*

  2. Never run a job with no workload remaining in the slice

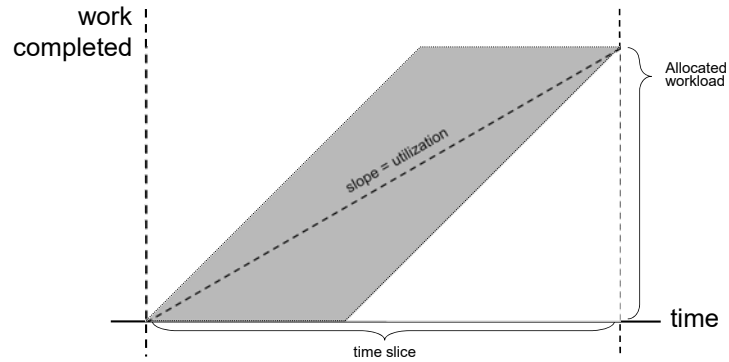  3. Do not voluntarily allow more idle processor time than $(m - \sum U_i) \times ($ *length of slice* $)$

## DP-Fair Work Allocation



work completed

slope = utilization

Allocated workload

time

time slice

## DP-Fair Scheduling Rule #2



work completed

When job finishes local workload, stop running

local work complete line

time

time slice

## DP-Fair Scheduling Rule #1



work completed

When job hits zero local laxity, then run to completion

zero local laxity line

time

time slice

## DP-Fair Scheduling Rule #3



idle time

Do not **voluntarily** allow idle time in excess of this limit

$slope = m$

$slope = m - \sum u_i$

Allowable idle time

time

time slice

## DF-Fair Guarantees Optimality

- We say that a scheduling algorithm is *DP-Fair* if it follows these three rules

- **Theorem:** Any DP-Fair scheduling algorithm for periodic tasks is optimal

---

## DP-Fair Implications

- ( Partition time into slices )
  + ( Assign proportional workloads )

  Optimal scheduling is almost trivial

  □ Minimally restrictive rules allow great latitude for algorithm design and adaptability

- What is the simplest possible algorithm?
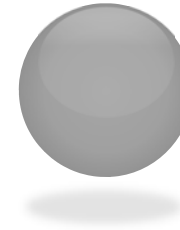
---

**EXAMPLE OF EXAM ASSIGNMENT: STUDYING THE RUN ALGORITHM**

PhD seminar on Real-Time Systems, University of Bologna, July 2014

---

RUN Assumptions

**Model parameters**

- $m$ homogeneous (symmetric) processors

- Implicit-deadline independent task $\tau_i, i \in \{1..n\}$

- $n = m + k, k \geq 0$

- Fixed-rate tasks $U_i = \frac{C_i}{T_i}$ $\qquad \sum_{i=1}^{n} U_i \leq m$

- Fully utilized system: no idle time (perhaps using fillers)

- *Migration* and *preemption* are assumed to have no additional costs over $c_i$

## Example /1



$n = 5$

**Legend**

task    Rate

processor

$k = n - m = 2$
(the excess)

$m = 3$

- $U_i = 0.6 \ \forall \tau_i, i = \{1, \dots, n = 5\}$
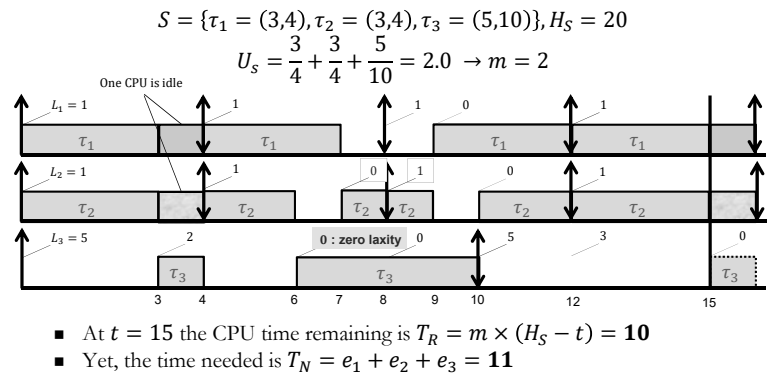- $\sum_{i=1}^{n=5} U_i = 3 \Rightarrow m = 3$ (fully utilized system)
- What schedule Σ for $\mathbf{S} = \{\{\tau_i\}, \mathbf{m}\}$ ?

## The G-LLF example at page 429 …

$$S = \{\tau_1 = (3,4), \tau_2 = (3,4), \tau_3 = (5,10)\}, H_S = 20$$
$$U_S = \frac{3}{4} + \frac{3}{4} + \frac{5}{10} = 2.0 \ \rightarrow m = 2$$



- At $t = 15$ the CPU time remaining is $T_R = m \times (H_S - t) = \mathbf{10}$
- Yet, the time needed is $T_N = e_1 + e_2 + e_3 = \mathbf{11}$

## Duality

- The problem of scheduling
  $\mathbf{S} = \{\tau_1 = (c_1, T_1), \dots, \tau_n = (c_n, T_n)\}, \mathbf{m}$
  has a *dual* problem that consists of scheduling
  $\mathbf{S'} = \{\tau'_1 = (T_1 - c_1, T_1), \dots, \tau'_n = (T_n - c_n, T_n)\}, (\mathbf{n} - \mathbf{m})$
- With this definition of duality
  - Laxity in primal is work remaining in the dual
  - A work-complete event in the primal is zero-laxity in the dual
  - And viceversa
- **Corollary**: any scheduling problem with $\mathbf{m}$ processors and $\mathbf{n} = \mathbf{m} + \mathbf{1}$ tasks and $\sum_1^n U_i = \mathbf{m}$ may be scheduled by applying EDF to its uniprocessor dual
  - If I can schedule $n$ tasks on $m$ processors, then I can also schedule the same $n$ tasks on $n - m$ processors
  - This is so because the scheduling events in the dual map to scheduling events in the primal

## Applying duality

$$S = \{\tau_1 = (3,4), \tau_2 = (3,4), \tau_3 = (5,10)\}, U_s = \frac{3}{4} + \frac{3}{4} + \frac{5}{10} = 2.0 \ \rightarrow m = 2$$

$$S_D = \{\tau_{1_D} = (1,4), \tau_{2_D} = (1,4), \tau_{3_D} = (5,10)\}, U_{S_D} = \frac{1}{4} + \frac{1}{4} + \frac{5}{10} = 1.0 \rightarrow m_D = 1$$



**The dual (LLF) schedule leaves no idle time**

Real-Time Systems     9

## Example /1

$n = 5$

**Legend**

task

Rate

processor

$k = n - m = 2$
(the excess)

$m = 3$

- $U_i = 0.6 \ \forall \tau_i, i = \{1, \dots, n = 5\}$
- $\sum_i^n U_i = 3 \Rightarrow m = 3$ (fully utilized system)
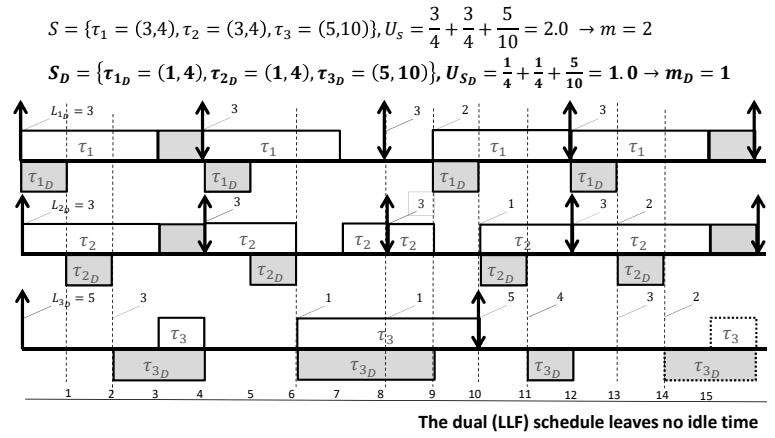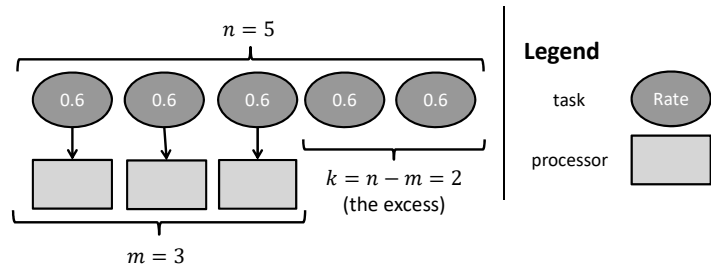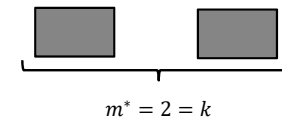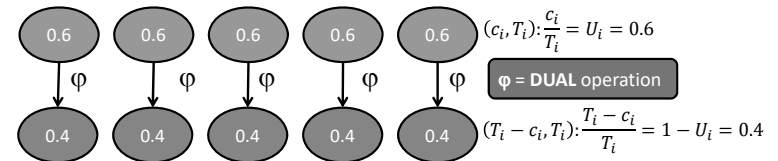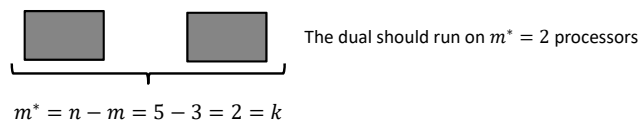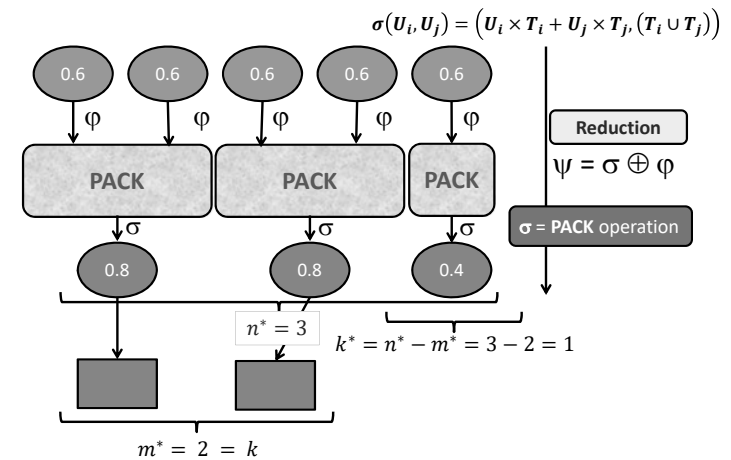- What schedule $\Sigma$ for $\mathbf{S} = \{\{\boldsymbol{\tau_i}\}, \boldsymbol{m}\}$ ?

## Example /3

0.6   0.6   0.6   0.6   0.6    $(c_i, T_i): \frac{c_i}{T_i} = U_i = 0.6$

$\varphi$   $\varphi$   $\varphi$   $\varphi$   $\varphi$    $\varphi$ = **DUAL** operation

0.4   0.4   0.4   0.4   0.4    $(T_i - c_i, T_i): \frac{T_i - c_i}{T_i} = 1 - U_i = 0.4$

$m^* = 2 = k$

## Example /2

- Consider the dual of this $\{n = 5, m = 3\}$ system

The dual should run on $m^* = 2$ processors

$m^* = n - m = 5 - 3 = 2 = k$

## Example /4

$\sigma(U_i, U_j) = \left( U_i \times T_i + U_j \times T_j, (T_i \cup T_j) \right)$

0.6   0.6   0.6   0.6   0.6

$\varphi$   $\varphi$   $\varphi$   $\varphi$   $\varphi$    **Reduction**

   $\psi = \sigma \oplus \varphi$

**PACK**   **PACK**   **PACK**

$\sigma$   $\sigma$   $\sigma$    $\sigma$ = **PACK** operation

0.8   0.8   0.4

$n^* = 3$   $k^* = n^* - m^* = 3 - 2 = 1$

$m^* = 2 = k$

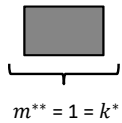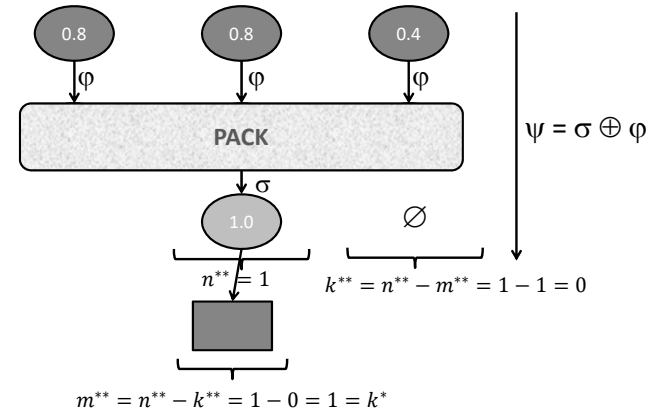## Example /5

The $(n^* = 3, m^* = 2)$ system still cannot be partitioned feasibly
Yet, applying duality to it seems promising since the dual would need $n^* - m^* = 1$ processor, which would REDUCE the problem TO a UNIPROCESSOR case
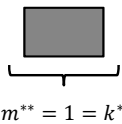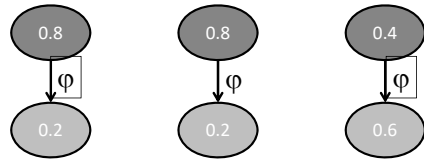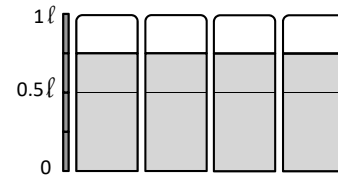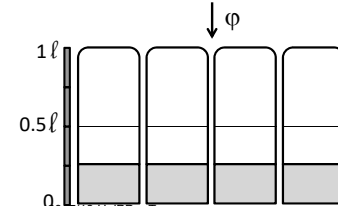
$$m^{**} = 1 = k^*$$

## Example /6



$$m^{**} = 1 = k^*$$

## Example /7



$$\psi = \sigma \oplus \varphi$$

$$n^{**} = 1 \qquad k^{**} = n^{**} - m^{**} = 1 - 1 = 0$$

$$m^{**} = n^{**} - k^{**} = 1 - 0 = 1 = k^*$$

## Why does reduction terminate? /1

**Lemma:** $\psi = \left| \sigma \circ \varphi \left( \bigcup_1^4 \tau_i \right) \right| \leq \left\lceil \dfrac{|\tau| + 1}{2} \right\rceil$



**Intuition**

$$\sum_1^4 U_i = 3 \Rightarrow m = 3$$
$$n = 4$$
$$\boldsymbol{k} = n - m = 1$$

**In the dual system**
$$\sum_1^4 U_i^* = n - m = 1 \Rightarrow$$
$$m^* = 1 = \boldsymbol{k}$$
$$n^* = 1 \text{ after packing}$$
$$k^* = 0 \text{ no leftover}$$

## Why does reduction terminate? /2

**Lemma:** $\psi = \left| \sigma \circ \varphi \left( \cup_1^4 \tau_i \right) \right| \leq \left\lceil \dfrac{|\tau|+1}{2} \right\rceil$



- Reduction $\psi = (\sigma \oplus \varphi)$ terminates as every step of it lowers the residual workload and the # of processors needed to run it
- The packing operation (at least) halves the number of tasks to schedule

- **Termination theorem:** after a finite number **$p$** of reduction steps, the system is reduced to a uniprocessor with full workload
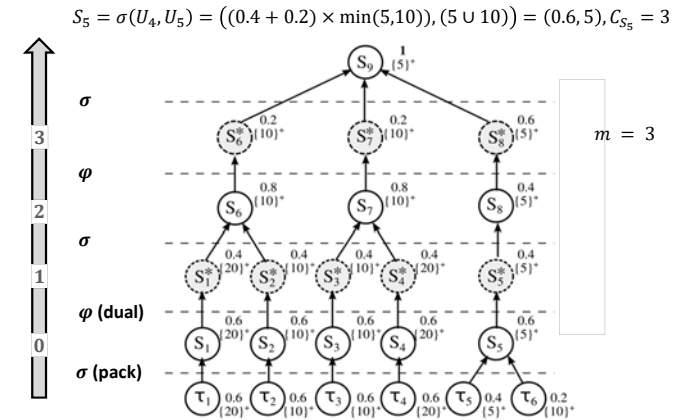
## How does RUN work /2

**Algorithm 1:** Outline of the RUN algorithm

I. OFF-LINE:
  A. Generate a reduction sequence for $\mathcal{T}$;
  B. Invert the sequence to form a server tree;
  C. For each proper subsystem $\mathcal{T}'$ of $\mathcal{T}$:
      Define the client/server at each virtual level;
II. ON-LINE:
Upon a scheduling event: ;
  A. If the event is a job release event at level $0$ ;
      1. Update deadline sets of servers on path up to root;
      2. Create jobs for each of these servers accordingly;
  B. Apply Rules 1 & 2 to schedule jobs from root to leaves, determining the $m$ jobs to schedule at level $0$;
  C. Assign the $m$ chosen jobs to processors, according to some task-to-processor assignment scheme;
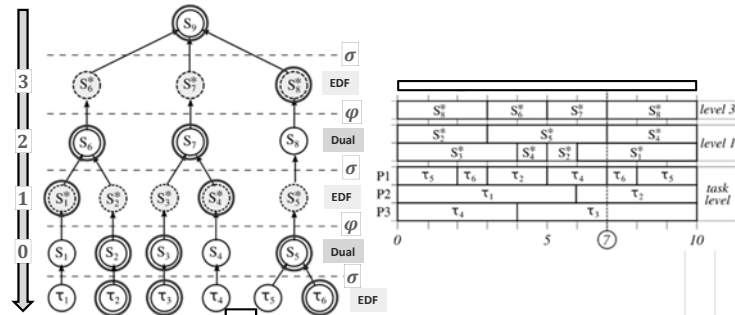
## How does RUN work /1

- A pair of basic operators
  - DUAL ($\varphi$)
  - PACK ($\sigma$)
- The REDUCE ($\psi = \sigma \oplus \varphi$) operation lowers (~ halves) the size of the problem at every step
- **Theorem** (validity of the dual): $\Sigma$ valid $\Leftrightarrow \Sigma^*$ valid
- Since every dual task represents the idle time of its primary, finding a feasible schedule for the dual (which is easier) determines a feasible schedule for its primary

## Example: off-line phase

$$S_5 = \sigma(U_4, U_5) = \left( (0.4 + 0.2) \times \min(5,10) \right), (5 \cup 10)) = (0.6, 5), C_{S_5} = 3$$

Example: on-line phase (at time $t = 7$)

## PROXIMA

# Putting RUN into practice
*Implementation and evaluation*

Davide Compagnin, Enrico Mezzetti and Tullio Vardanega
University of Padua, Italy

26[th] EUROMICRO Conference on Real-time Systems (ECRTS)
Madrid, 9 July 2014

*www.proxima-project.eu*

---

## RUN implementation

❑ **For real**
  ➢ On top of **LITMUS[RT]** Linux test-bed (UNC, now MP-SWI)
  ➢ Relying on *standard* RTOS support
❑ Main implementation choices and challenges
  ➢ *Scheduling on the reduction tree*
    - How to organize the data structure
    - How to perform virtual scheduling and trigger tree updates
    - Intrinsic influence of the packing policy
  ➢ *Mixing global and local scheduling*
    - Global release event queue vs. local *level-0* ready queue
    - Handling simultaneous scheduling events
      • Job release, budget exhaustion (possibly from different sub-trees)
  ➢ *Meeting the full-utilization requirement*
    - Variability of tasks' WCET and less-than-full utilization

---

## Empirical evaluation

❑ **Empirical evaluation** instead of simulation-based

❑ Focus on **scheduling interference**
  ➢ Cost of scheduling primitives
  ➢ Incurred preemptions and migrations

❑ RUN compared against **P-EDF** and **G-EDF**
  ➢ RUN shares something in common with both
  ➢ Much better than **Pfair** (S-PD$^2$ in LITMUS[RT])
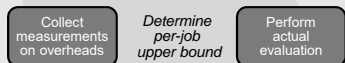    - RUN has superior performance for preemptions and migrations

## Experimental setup

- ❑ **LITMUS^RT** on an 8-core AMD Opteron™ 2356

- ❑ Collected measurements for RUN, P-EDF, G-EDF
  - ➢ Hundreds of automatically generated task sets
  - ➢ Harmonic and non-harmonic, with global utilization @ 50%-100%
  - ➢ Representative of small up to large tasks

- ❑ **Two-step process**
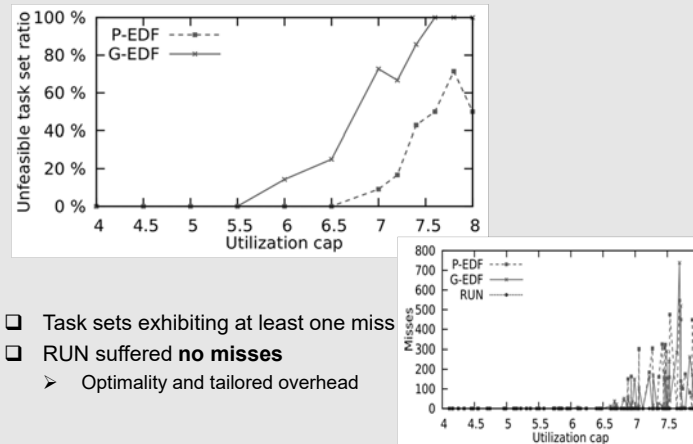  - ➢ Preliminary empirical determination of overheads

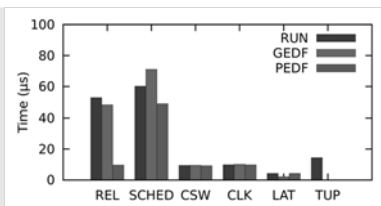505                                                  09/07/2014   **PROXIMA**

## Empirical schedulability

- ❑ Task sets exhibiting at least one miss
- ❑ RUN suffered **no misses**
  - ➢ Optimality and tailored overhead

507                                                  09/07/2014   **PROXIMA**

## Primitive overheads and empirical bound

- ❑ Expectations confirmed
  - ➢ P-EDF needs lighter-weight scheduling primitives
- ❑ **Tree update** (TUP) triggered upon
  - ➢ *Budget exhaustion* event
  - ➢ Job release → REL includes TUP
- ❑ Empirical upper bound on RUN scheduling overhead
  - ➢ $OH_{RUN}^{Job} = REL + \widehat{SCHED} + CLK + k \times (TUP + \widehat{SCHED} + max(PRE, MIG))$

$$k = \lceil (3p+1)/2 \rceil \quad \text{and} \quad \widehat{SCHED} = SCHED + CSW + LAT.$$

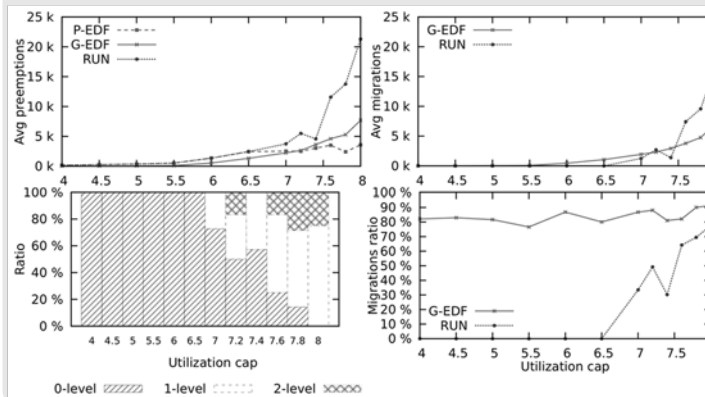506                                                  09/07/2014   **PROXIMA**

## Kernel interference

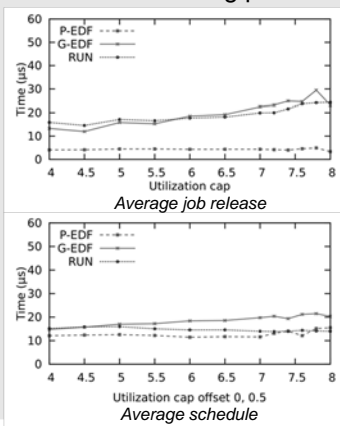- ❑ Observing average preemptions and migrations

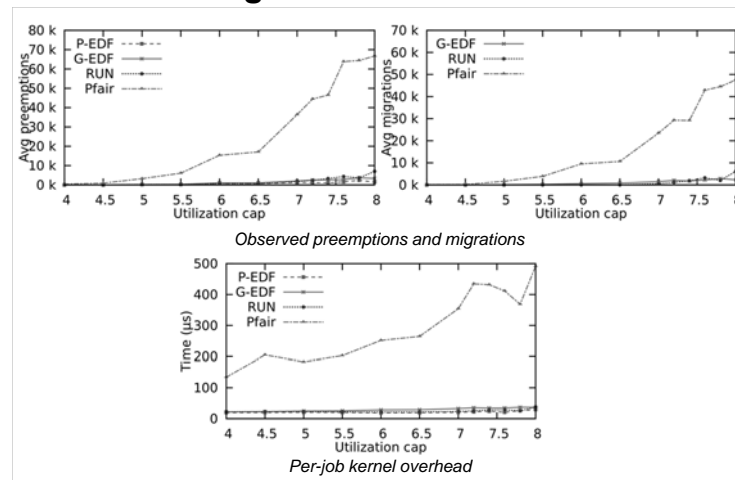508                                                  09/07/2014   **PROXIMA**

## Scheduling cost

❑ Average cost of core scheduling primitives



*Average job release*

*Average schedule*

509 09/07/2014 **PROXIMA**

## Evaluation against S-PD²



*Observed preemptions and migrations*

*Per-job kernel overhead*

511 09/07/2014 **PROXIMA**

## Per-job scheduling overhead



*Harmonic task set*        *Non-harmonic task set*

510 09/07/2014 **PROXIMA**