



Global Scheduling on Heterogenous MPSoCs

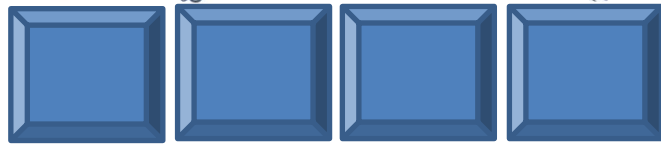
Emmanuel Grolleau
grolleau@ensma.fr

Based on a work done with :

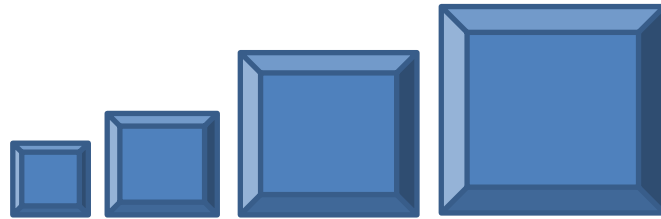
Antoine Bertout, Joël Goossens, Xavier Poczekajlo, Roy Jamil

Classification of mp platforms

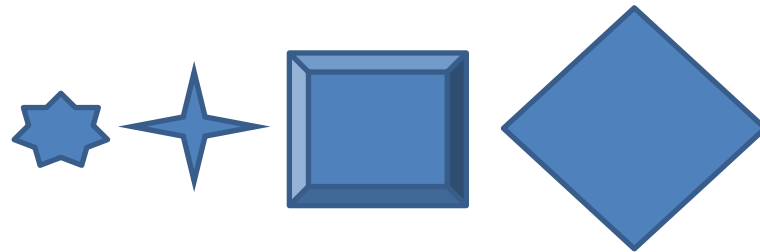
Identical



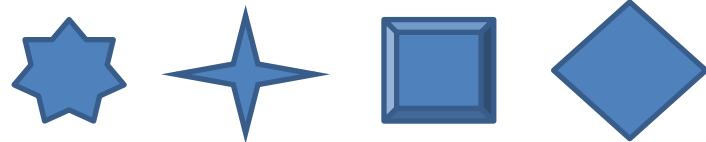
Uniform



Consistent



Unrelated



	Π_1	Π_2	Π_3	Π_4
τ_i	Orange square	Orange square	Orange square	Orange square
τ_j	Green square	Green square	Green square	Green square
τ_i	Orange square	Orange square	Orange square	Orange rectangle
τ_j	Green square	Green square	Green square	Green rectangle
τ_i	Orange square	Orange square	Orange square	Orange rectangle
τ_j	Green square	Green square	Green rectangle	∞
τ_i	Orange square	Orange square	Orange square	Orange square
τ_j	Green square	Green square	∞	Green rectangle

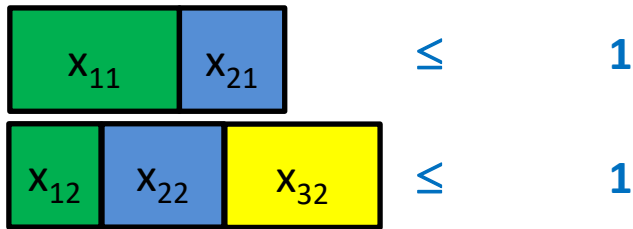
- Seminal paper in operations research scheduling Lawler & Labetoulle 1978
- Seminal paper in r-t scheduling of S. Baruah 2004
 - Input: (strictly) periodic independent synchronous implicit deadline tasks, defined by $u_i=C_i/T_i$, C_i defined on a fictional reference core, and for each core Π_j , a rate r_{ij}

utilization	$\Pi_{\text{fictional}}$
τ_1	u_1
τ_2	u_2
τ_3	u_3

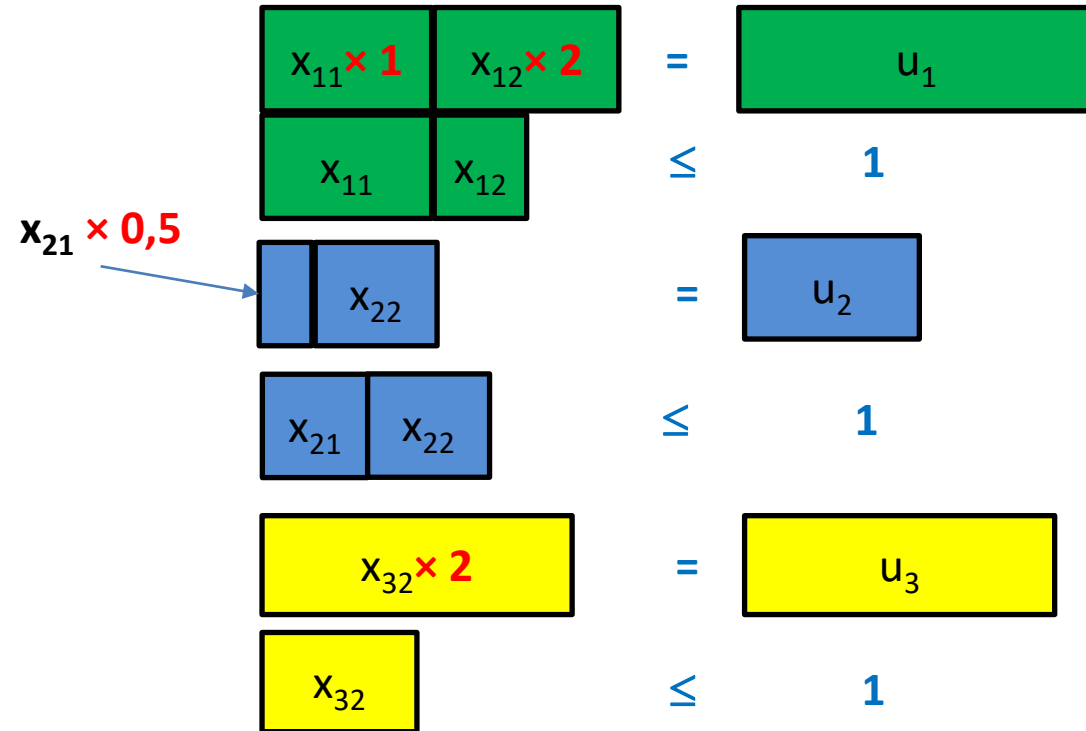
rates	Π_1	Π_2
τ_1	$r_{11}=1$	$r_{12}=2$
τ_2	$r_{21}=0.5$	$r_{22}=1$
τ_3	$r_{31}=0$	$r_{32}=2$

- Linear Program: what fraction of core to what task?
- Theorem: the system is feasible **if and only** if the LP has a solution

Workload assignment	Π_1	Π_2
τ_1	$x_{11} \times 1$	$x_{12} \times 2$
τ_2	$0,5 \times x_{21}$	$x_{22} \times 1$
τ_3		$2 \times x_{32}$

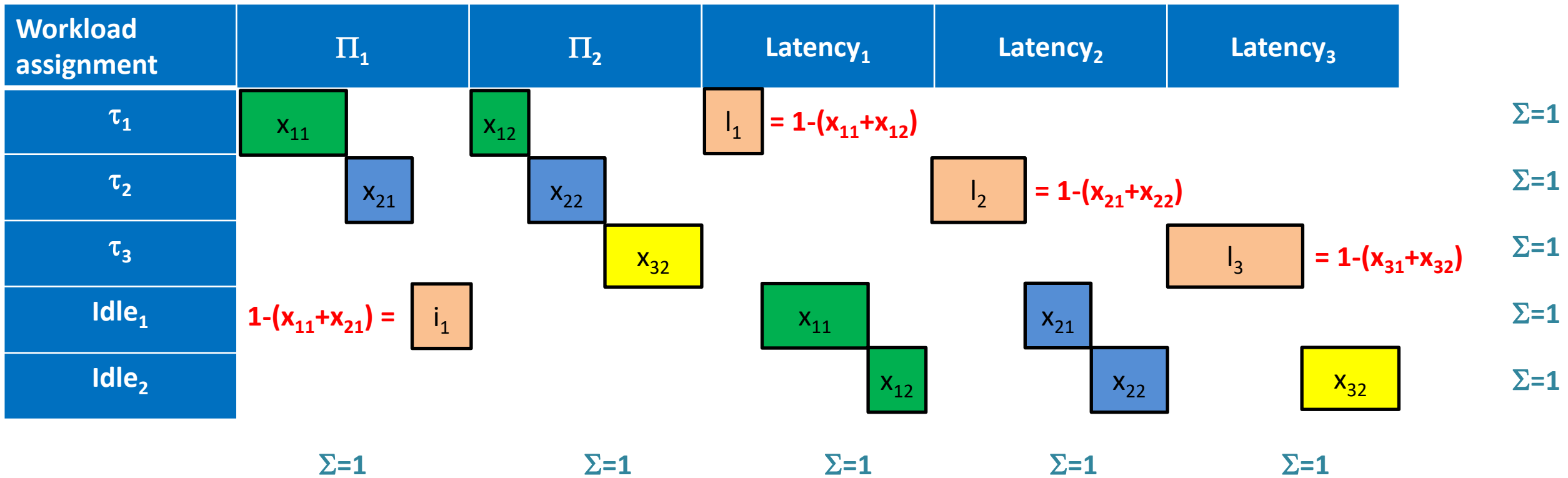


$$\forall \text{core } j, \sum_{\text{task } i} x_{ij} \leq 1$$



$$\forall \text{task } i, \sum_{\text{core } j} x_{ij} r_{ij} = u_i$$

$$\forall \text{task } i, \sum_{\text{core } j} x_{ij} \leq 1$$



□ This is a Doubly Stochastic (DS) matrix!

□ A DS matrix

➤ Square, non negative values, sum of each row and column is 1

□ A DS matrix can be expressed as a convex combination of permutation matrices

$$\begin{array}{|c|c|c|} \hline 0.3 & 0.4 & 0.3 \\ \hline 0.5 & 0.5 & 0 \\ \hline 0.2 & 0.1 & 0.7 \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 0 & 0 & 1 \\ \hline 1 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline \end{array} \times 0.1 + \begin{array}{|c|c|c|} \hline 0.3 & 0.4 & 0.2 \\ \hline 0.4 & 0.5 & 0 \\ \hline 0.2 & 0 & 0.7 \\ \hline \end{array}$$

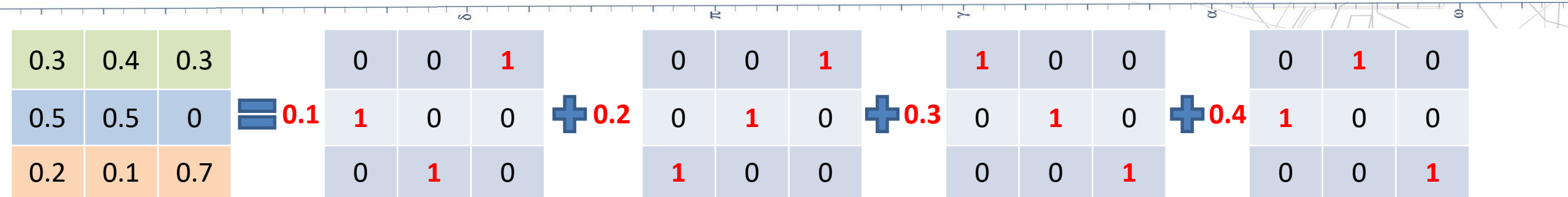
□ Note that $1/0.9 \times$

0.3	0.4	0.2
0.4	0.5	0
0.2	0	0.7

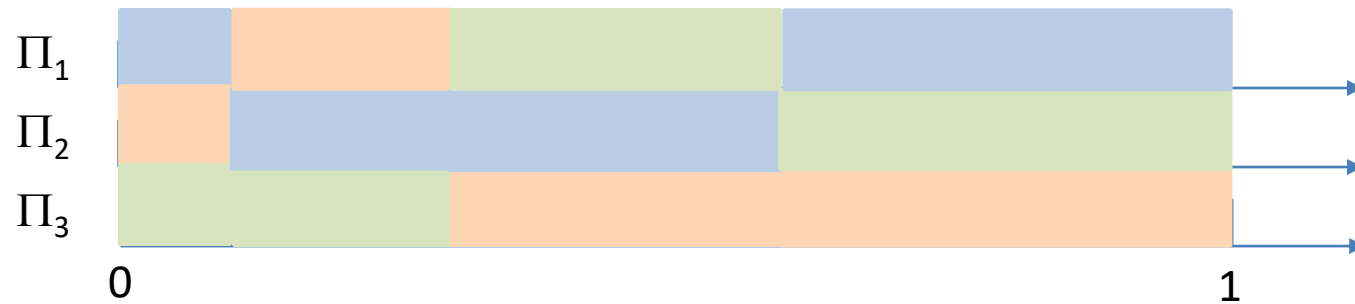
is DS

Convex combination of permutation matrices

$$\begin{bmatrix} 0.3 & 0.4 & 0.3 \\ 0.5 & 0.5 & 0 \\ 0.2 & 0.1 & 0.7 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \times 0.1 + \begin{bmatrix} 0.3 & 0.4 & 0.2 \\ 0.4 & 0.5 & 0 \\ 0.2 & 0 & 0.7 \end{bmatrix} \\
 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \times 0.2 + \begin{bmatrix} 0.3 & 0.4 & 0 \\ 0.4 & 0.3 & 0 \\ 0 & 0 & 0.7 \end{bmatrix} \\
 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \times 0.3 + \begin{bmatrix} 0 & 0.4 & 0 \\ 0.4 & 0 & 0 \\ 0 & 0 & 0.4 \end{bmatrix}$$



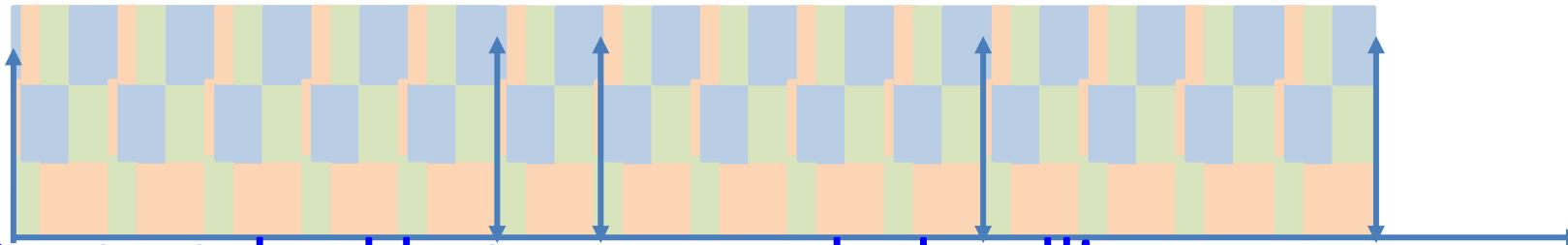
□ As a real-time scheduling guy, I can interpret it as a schedule...



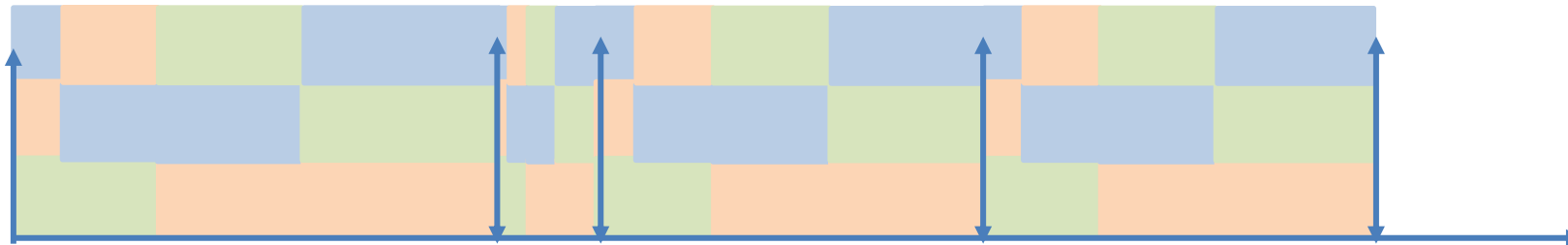
□ Each permutation matrix generates a scheduling point

➤ Preemption and/or migration

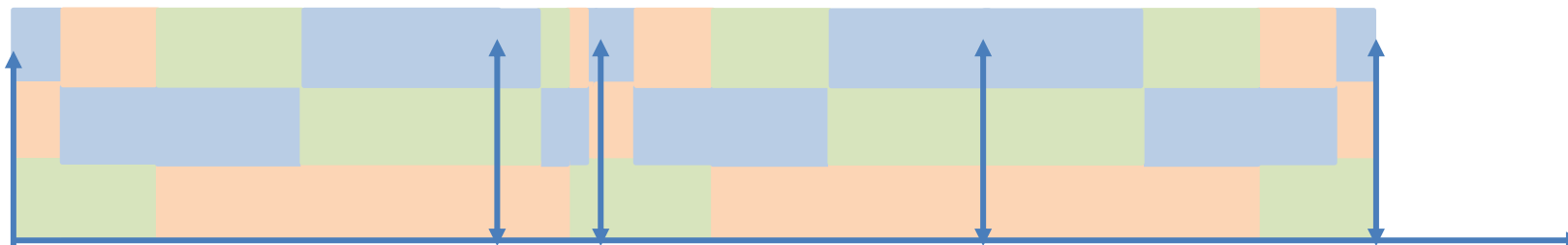
□ A template schedule can be repeated on each time unit



□ Or stretched between each deadline

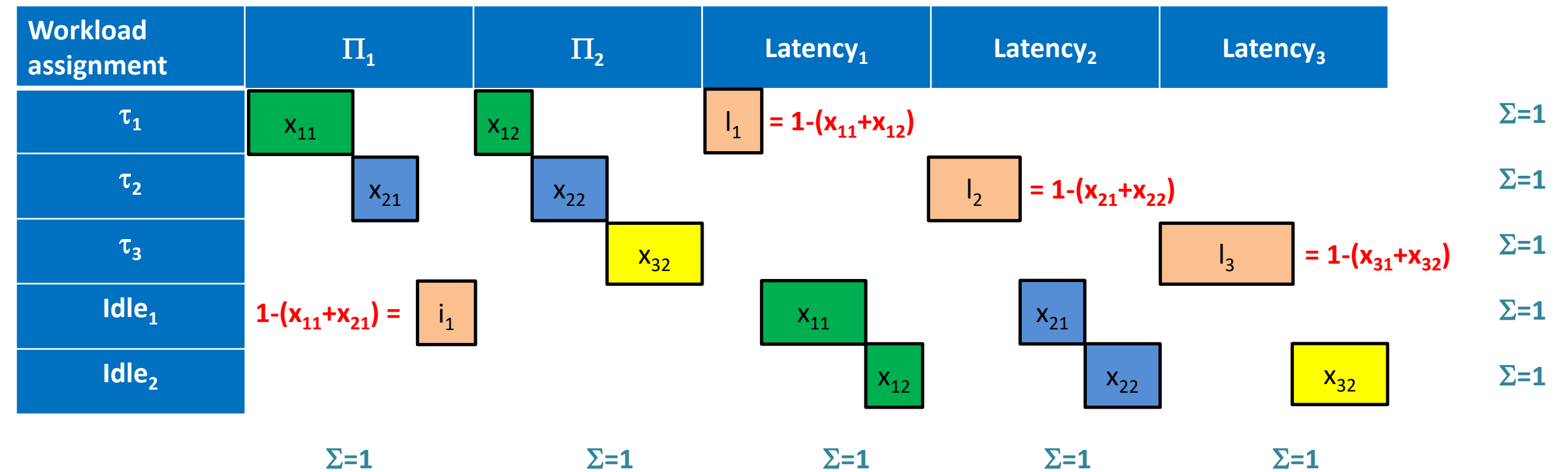


□ Mirrored every other time

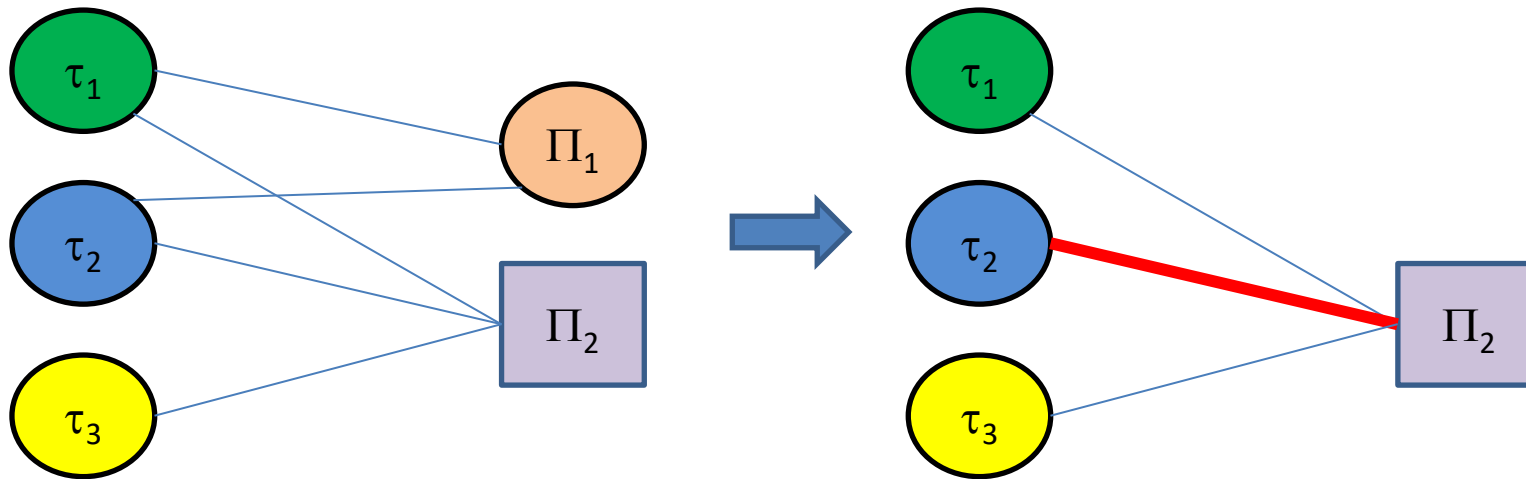


- ❑ Can always be obtained from a DS matrix (BvN decomposition theorem)
- ❑ Obtaining a valid workload assignment matrix is a necessary and sufficient schedulability condition
 - Under the hypothesis of no preemption cost, no migration cost
- ❑ The produced off-line schedule supposes an « almost fluid » scheduler, able to preempt/migrate tasks several times per time unit
- ❑ Number of permutation matrices = number of scheduling points per template schedule

S. Baruah's strategy

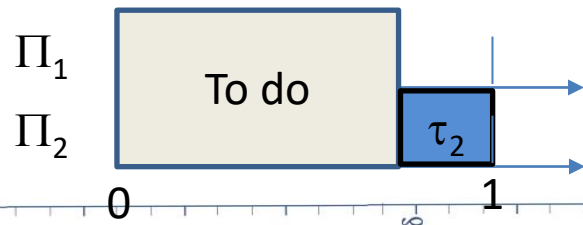
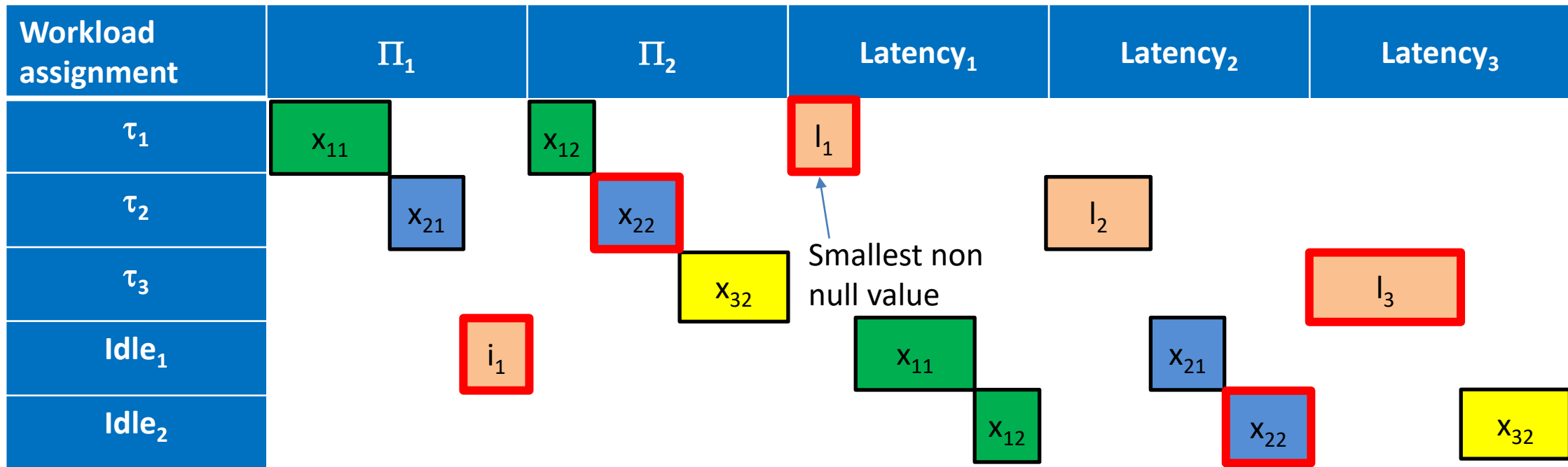


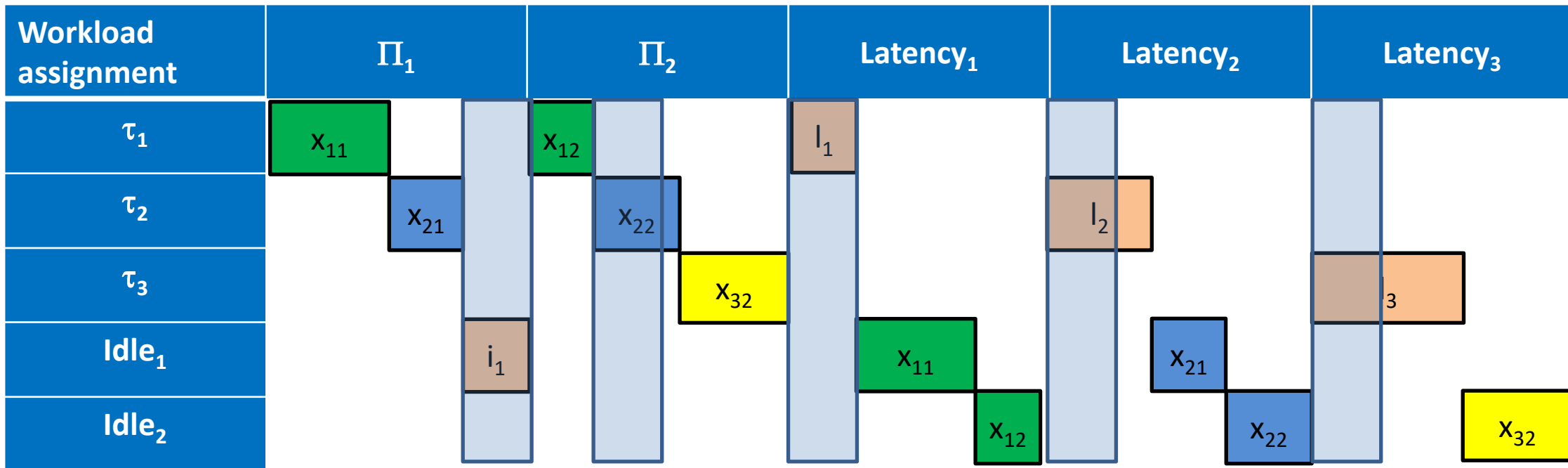
- Square node = urgent task or full processor, assign absolutely
- Circle node = non urgent, non full, assign if necessary



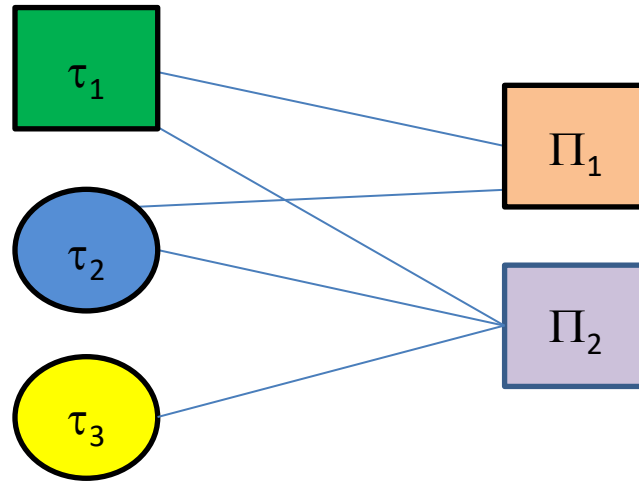
Arbitrary marriage

Reverse construction of template schedule





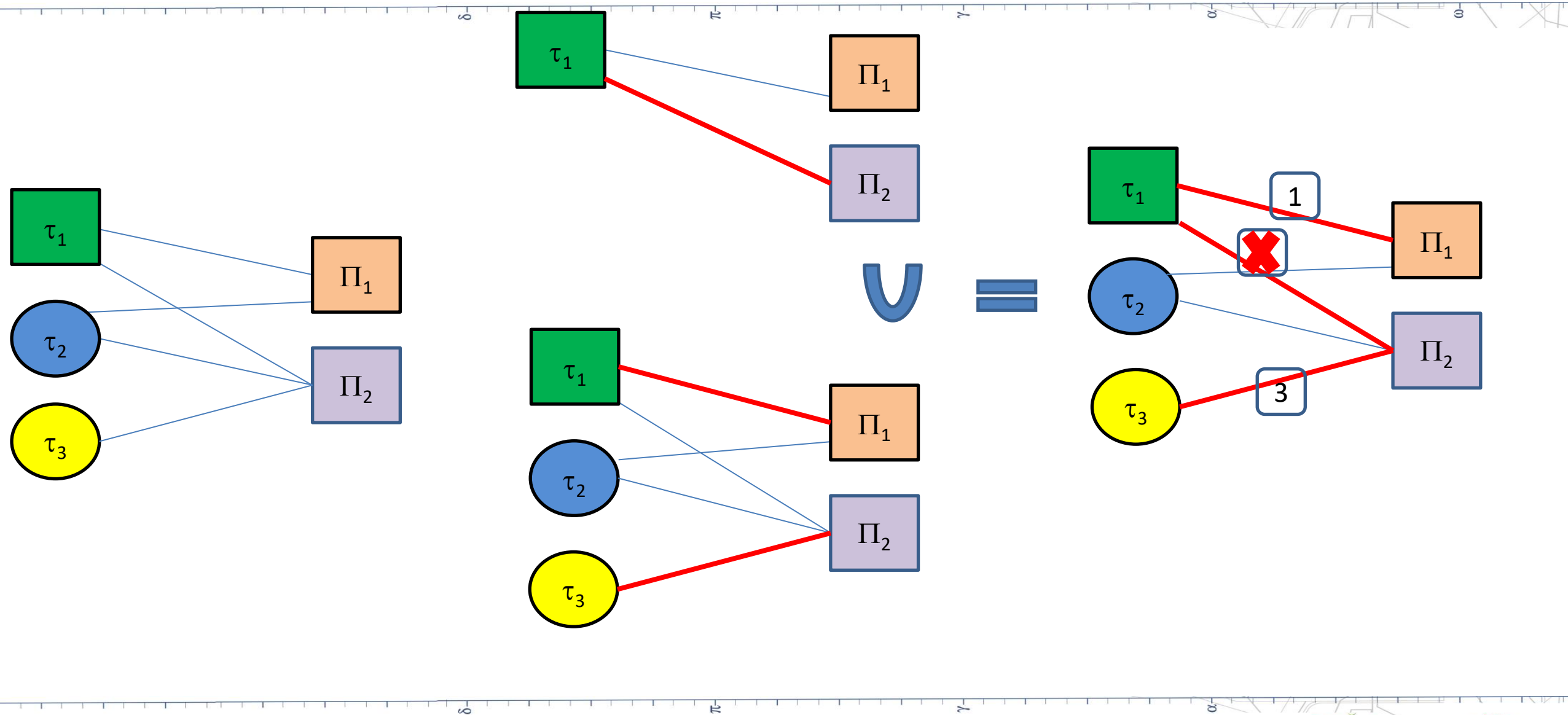
- Π_1 is now full, like Π_2 , and they will remain full until the... beginning of the template schedule
- τ_1 is now urgent, and will remain urgent until the beginning

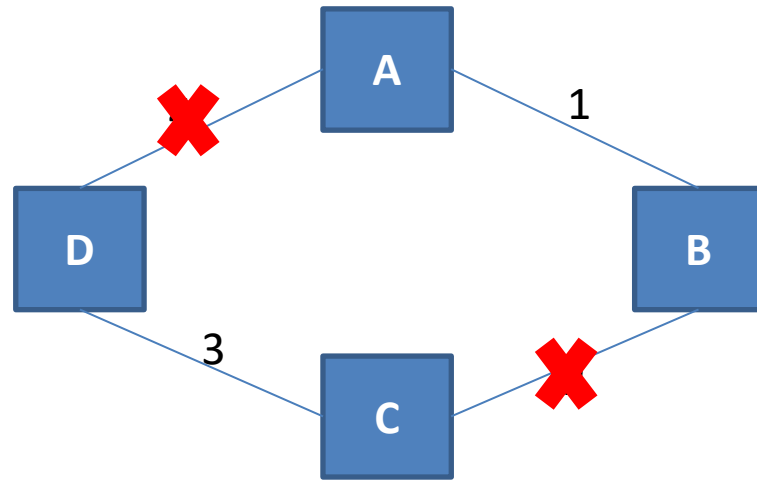
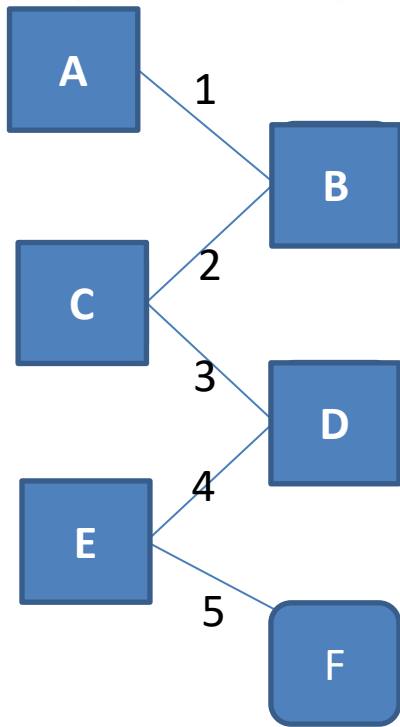


- Find a marriage such that each square node is married
- Circle nodes are "spare nodes"

- Could we invent a new "marriage in the nobility" problem?

Baruah's method

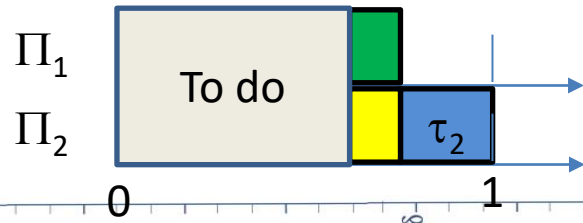
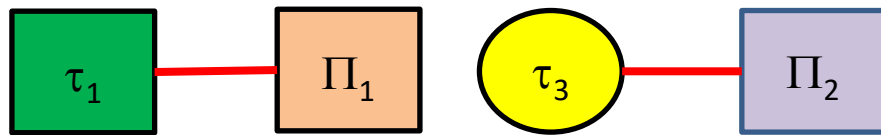
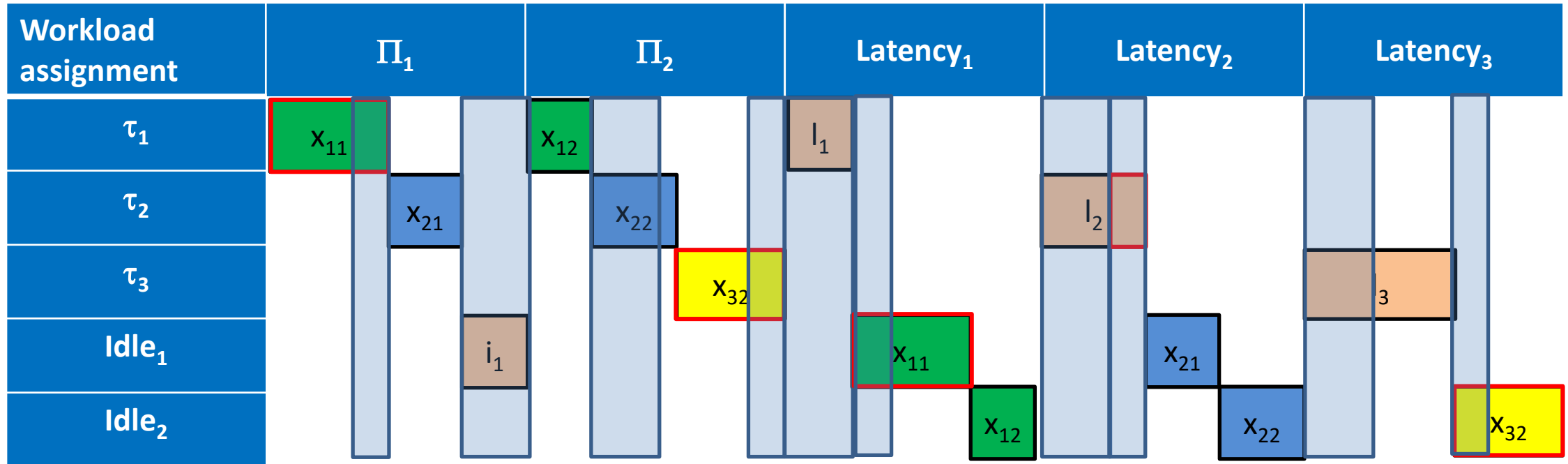




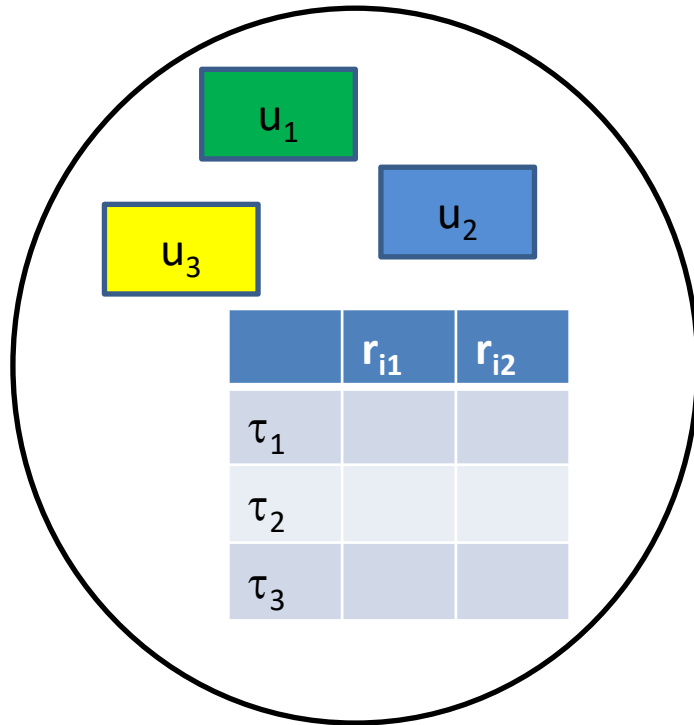
- Number any path starting from an important node
- Cut the edges with an even number

- Only important nodes can be the endpoints of two edges
- An non cyclic even length path has a non important node as one of its extremities
- A cycle can only have an even length path

Template schedule construction



Tasks + rates on cores



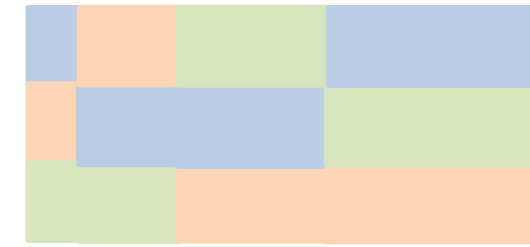
LP

Workload assignment

	Π_1	Π_2
τ_1	x_{11}	x_{12}
τ_2	x_{21}	x_{22}
τ_3	x_{31}	x_{32}

BvN

Template schedule



$$\forall \text{task } i, \sum_{\forall \text{core } j} x_{ij} r_{ij} = u_i$$

$$\forall \text{task } i, \sum_{\forall \text{core } j} x_{ij} \leq 1$$

$$\forall \text{core } j, \sum_{\forall \text{task } i} x_{ij} \leq 1$$

$$\forall \text{task } i, \sum_{\forall \text{core } j} x_{ij} r_{ij} = u_i$$

$$\forall \text{task } i, \sum_{\forall \text{core } j} x_{ij} \leq L$$

$$\forall \text{core } j, \sum_{\forall \text{task } i} x_{ij} \leq L$$

LP Feas Obj: min L

□ The objective function can be any

LP Load Obj: min $\sum_j \sum_i x_{ij}$

□ Room left for e.g. energy or heat dissipation optimization

□ Constraints can be added

$$\forall \text{task } i, \forall \text{core } j, b_{ij} \in \{0,1\}$$

$$x_{ij} \leq b_{ij}$$

$$b_{ij} < 1 + x_{ij}$$

} b_{ij} is 1 iff τ_i uses Π_j , 0 else

ILP Mig Obj: min $\sum_j \sum_i b_{ij}$

LP Feas vs. LP Load

	C_i	T_i
τ_1	5	10
τ_2	5	10

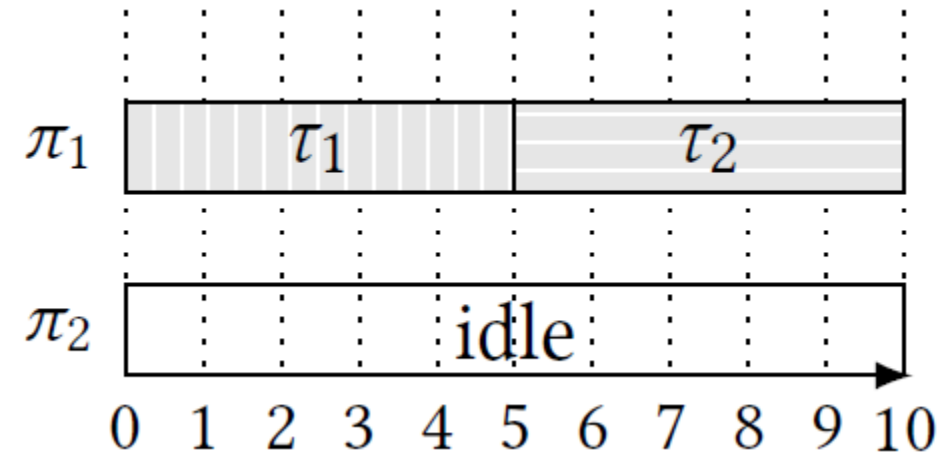
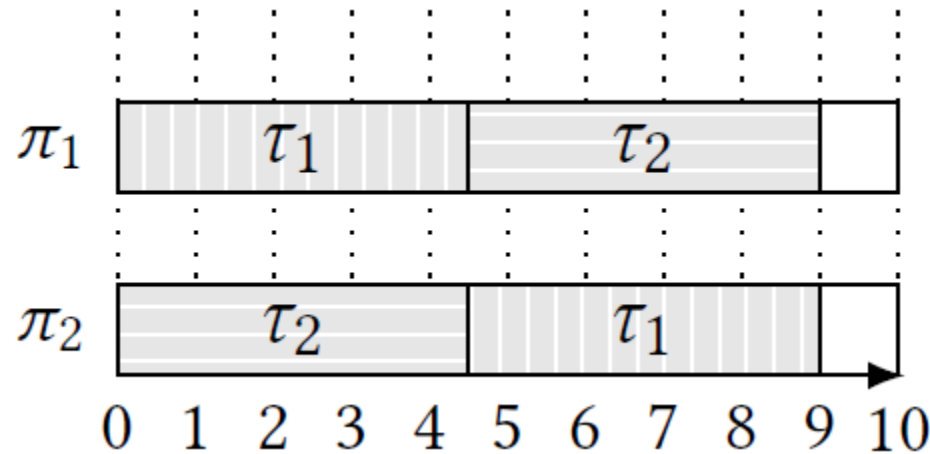
Rates	π_1	π_2
τ_1	10	1
τ_2	10	1

LP Feas

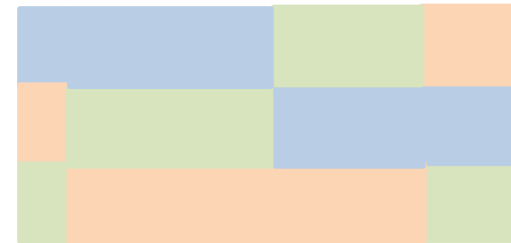
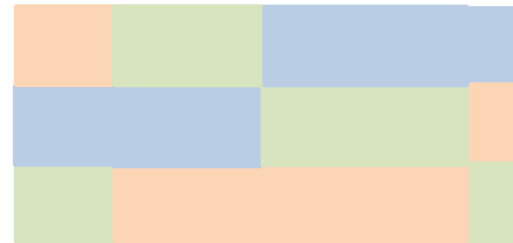
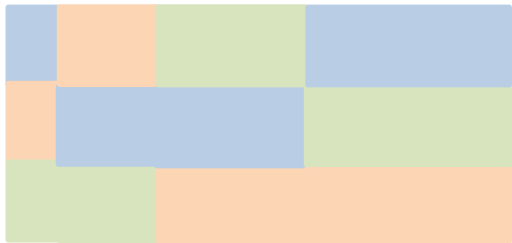
LP Load

$$\begin{bmatrix} 5/11 & 5/11 \\ 5/11 & 5/11 \end{bmatrix}$$

$$\begin{bmatrix} 0,5 & 0,5 \\ 0 & 0 \end{bmatrix}$$



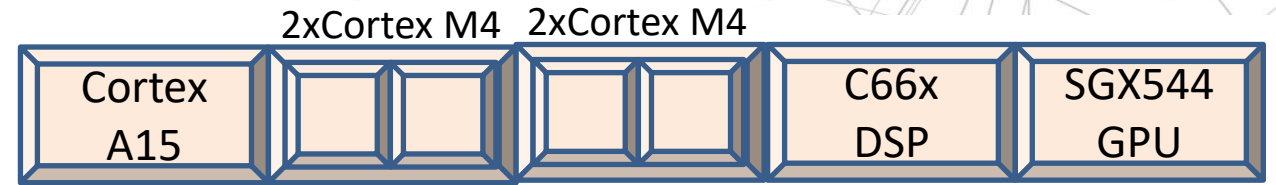
- ❑ Metrics: number of migrations & preemptions
- ❑ And (experimentally) the winner is...
- ❑ The conservative decomposition



- ❑ Traveler salesman problem

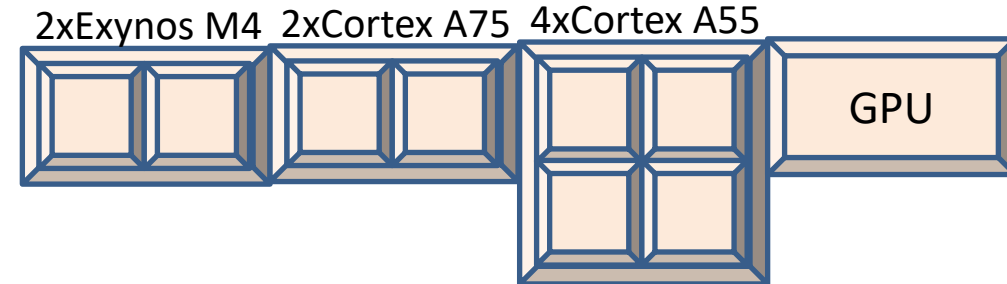
TI Sitara AM57x

➤ Embedded computing, robotics, avionics, medical imaging, etc.



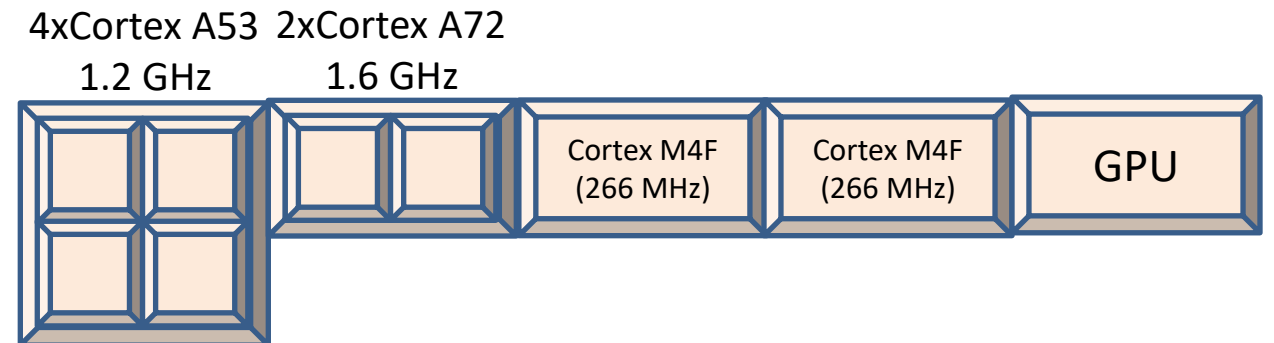
Samsung Exynos 9 9820

➤ Smartphone

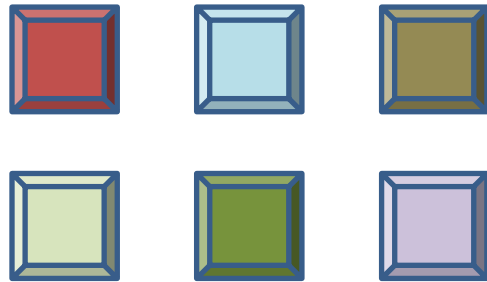


NXP i.MX 8 QuadMax

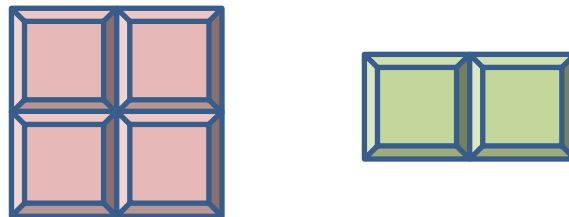
➤ Automotive, etc.



❑ Rather than having a heterogeneous platform



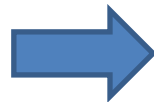
❑ Consider a set of clusters of identical cores



$$\forall \text{task } i, \sum_{\forall \text{cluster } j} x_{ij} r_{ij} = u_i$$

$$\forall \text{task } i, \sum_{\forall \text{cluster } j} x_{ij} \leq 1$$

$$\forall \text{cluster } j, \sum_{\forall \text{task } i} x_{ij} \leq m_j$$

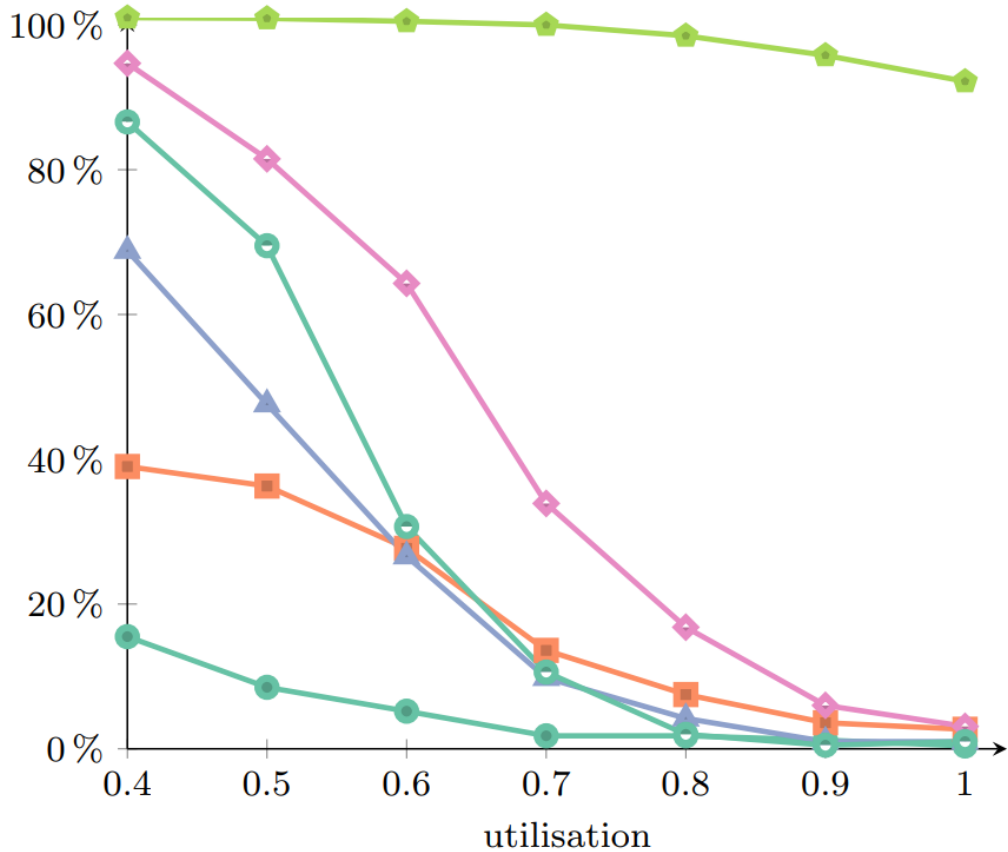


Workload assignment	Π_1	Π_2	
τ_1	x_{11}	x_{12}	$\Sigma \leq 1$
τ_2	$0,5 \times$	$1 \times$	$\Sigma \leq 1$
τ_3	$2 \times$		$\Sigma \leq 1$
	$\Sigma \leq m_1$	$\Sigma \leq m_2$	

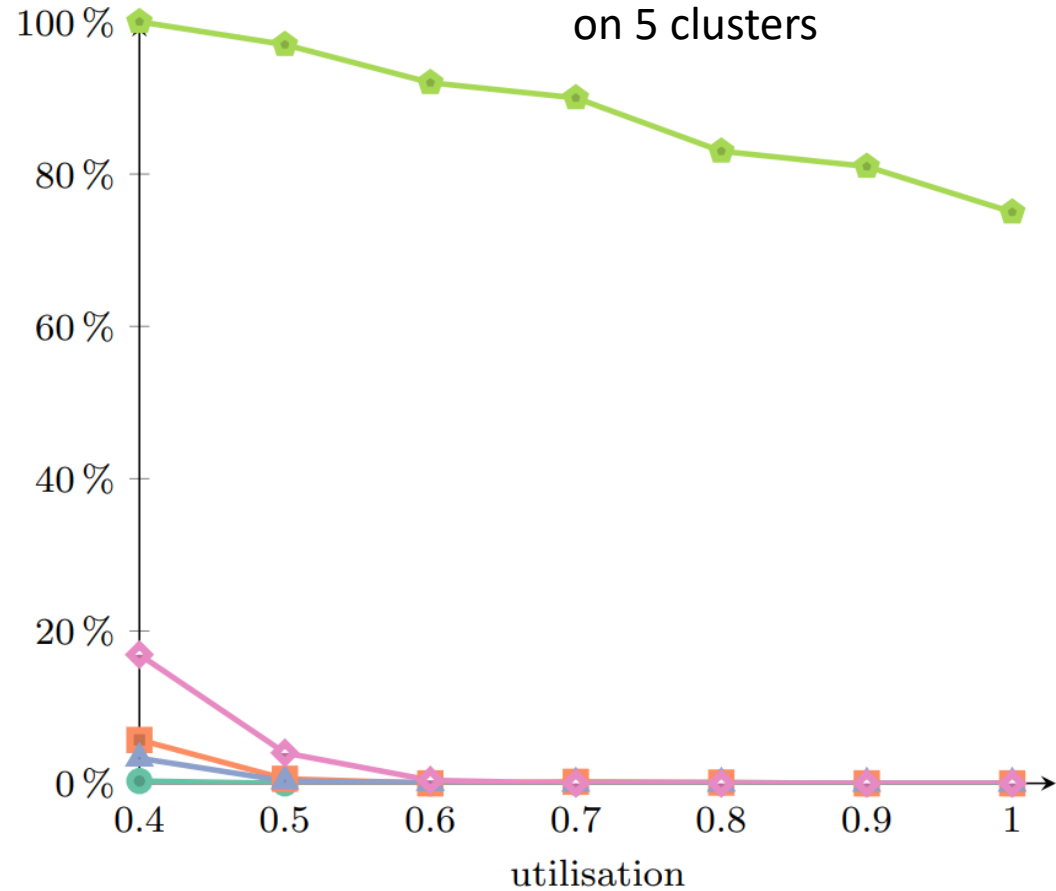
- ❑ A system is feasible iff LP has a solution => always possible to build a DS matrix
- ❑ Less variables (rate r_{ij} per cluster)
 - ILP for inter-cluster migrations minimization smaller
- ❑ Inter-cluster \neq Intra-cluster migration
 - Experimentally 10 to 70 μ s vs. 1 to 2 μ s on i.mx8 and STM32MP1

Comparison flat vs. clustered

Average number of clustered workload assignments on 2 clusters



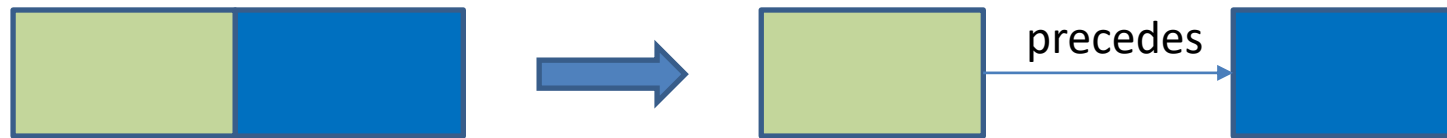
Average number of clustered workload assignments on 5 clusters



□ Average execution time (in seconds)

	2 clusters	5 clusters
LP-Feas	0.013	0.464
LP-Load	0.012	0.562
LP-CFeas	0.002	0.027
LP-CLoad	0.002	0.029
Hetero-split	0.007	N/A
ILP-Mig	0.061	N/A
ILP-CMig	0.023	0.156

- ❑ Zero cost for preemption & migration => already NP-hard for uniprocessor
- ❑ « If any portion of a thread is executed for 5% of the time on a core of rate 2, it executes for 10% »
 - What if the first half of a thread uses intensively integers, and the second half uses intensively floats?



- ❑ Limited to implicit deadlines, strictly periodic tasks
 - Sporadic tasks with explicit deadlines?
- ❑ Offline schedule hard to implement
 - Dynamic schedule à la U-EDF?

- ❑ Are used and will be more in the future
- ❑ Energy saving possibilities
 - DVFS, DPM
- ❑ Global scheduling
 - Can use up to 100% of the platform
 - Can be seen as saving more energy in the future
- ❑ Lots to do...