

PARTITION WITH SIDE EFFECTS

Krzysztof (Chris) Rządca

Institute of Informatics, University of Warsaw, Poland

joint work with Fanny Pascual

Sorbonne Universités, Université Pierre et Marie Curie, LIP6, Paris

CLOUD COMPUTING: THE FRONT- END OF THE MODERN DATACENTER



Google Cloud Platform Live

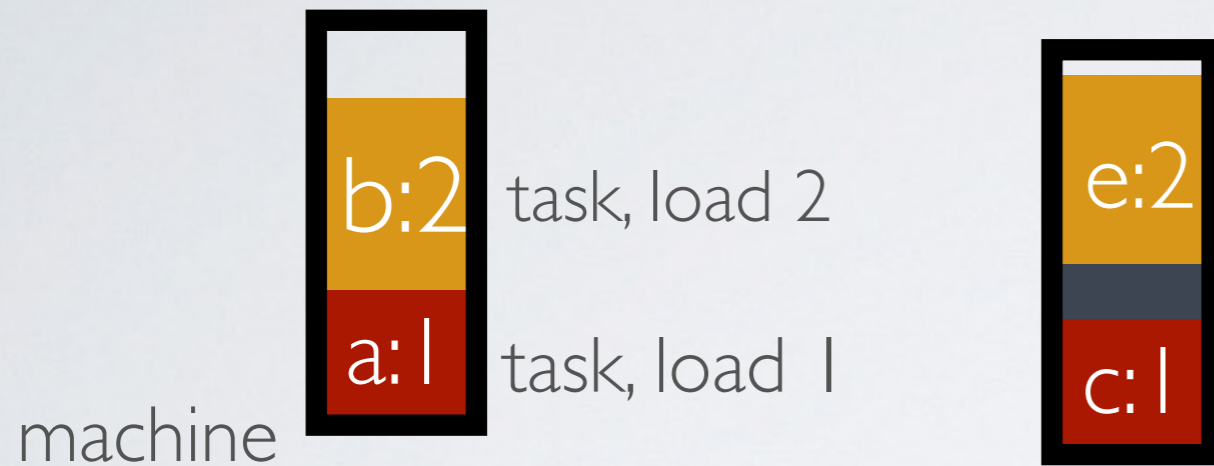
- virtual machines for hire for 0.10\$-2.00\$ per hour
- used by many organizations to reduce infrastructure costs

A DATACENTER IS NOT YOUR HPC SUPERCOMPUTER

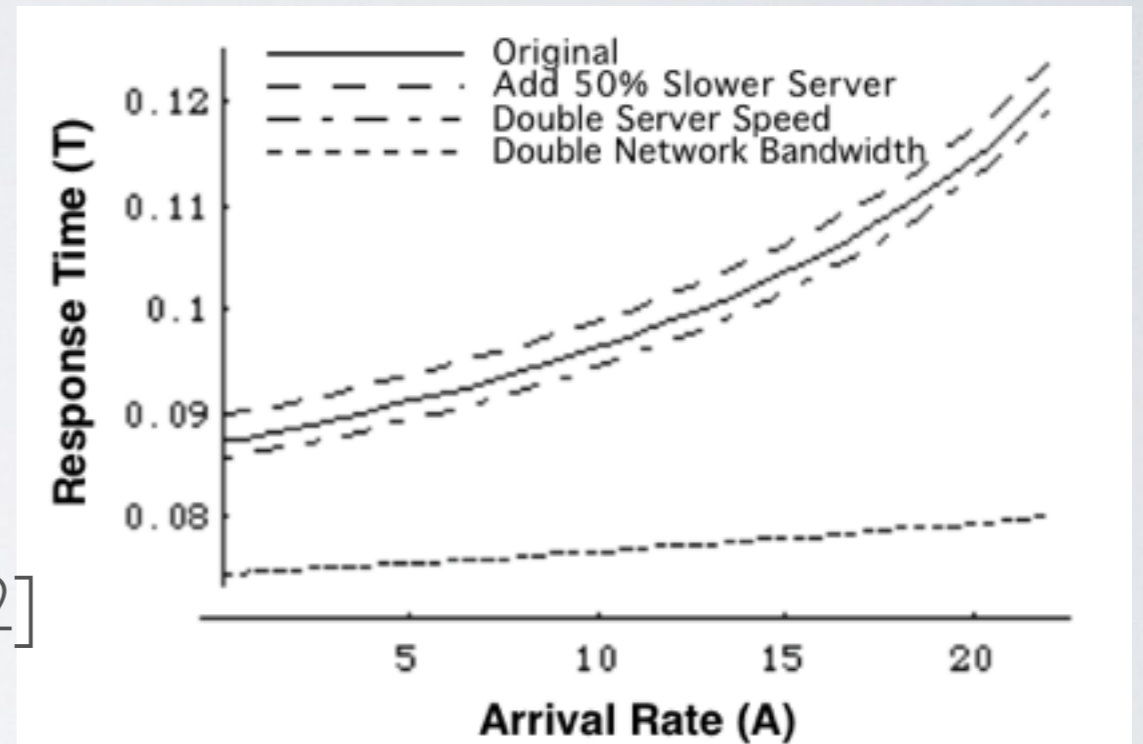
- **co-location of dozens of tasks on a single physical node (focus of this talk)**
- virtualization&migrations
- Service Level Agreements
- surviving failures
- considering network bandwidth
-

CO-ALLOCATING TASKS ONTO MACHINES

bin packing?



optimizing user experience?



multi-dimensional bin packing?

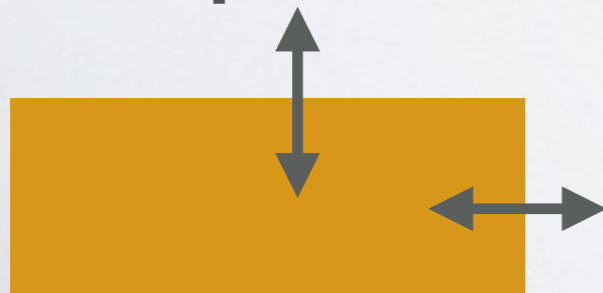
[Stillwell et al 2012]



[Slothouber 95]

loads are probabilistic?

[Goel&Indyk99, WMZII, ...]



CONSTRUCTING THE MODEL: TASKS ARE HETEROGENEOUS AND IMPACT EACH OTHER IN A DIFFERENT WAY!

consider a web service with load 2



and a RAM-cached database with load 1



They use different resources
(webservice: mostly network; db: mostly RAM/cache)

colocating a web service
with a database



>
(performance)



colocating
2 web services

OUR PERFORMANCE MODEL: A TASK HAS **LOAD** AND **TYPE**

a:2 load=2
type=webservice

b:1 load=1
type=db

c:3 load=3
type=db

task's cost = f_t (sum of loads of type web; sum of loads of type db)
(cost = -performance)

a:1
b:1
d:3

$c_a = f_{db}(3; 2)$

M1

d:1
e:3

$c_e = f_{web}(3; 1)$

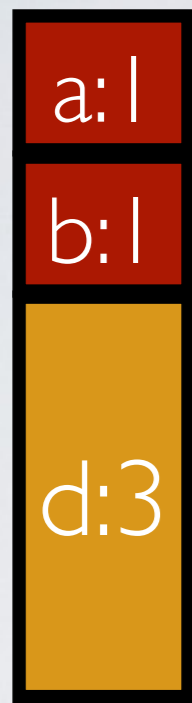
M2

goal: $\min \sum \text{cost}$

$$\sum \text{cost} = 2f_{db}(3; 2) + f_{db}(3; 1) + f_{web}(3; 2) + f_{web}(3; 1)$$

LINEAR PERFORMANCE MODEL:

$$\text{COST} = \text{COEFFICIENT} * \text{LOAD}$$



$$c_a = f_{db}(3; 2) = w[\text{web}, \text{db}] * 3 + w[\text{db}, \text{db}] * 2$$



$$c_e = f_{web}(3; 1) = w[\text{web}, \text{web}] * 3 + w[\text{db}, \text{web}] * 1$$

M1

M2

$$\sum \text{cost} = 2f_{db}(3; 2) + f_{db}(3; 1) + f_{web}(3; 2) + f_{web}(3; 1)$$

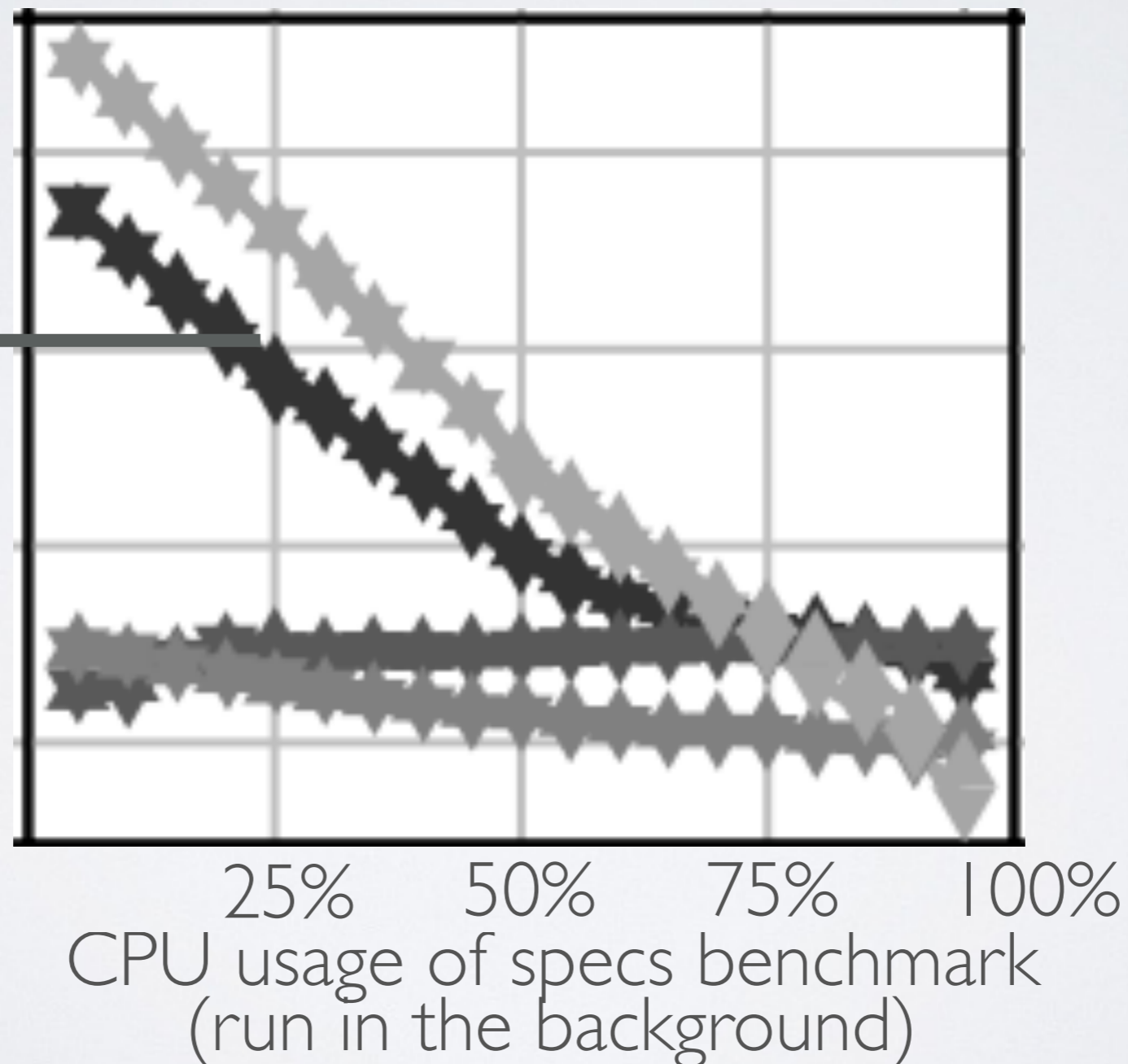
$$\begin{aligned} \sum \text{cost} = & 2 * (w[\text{web}, \text{db}] * 3 + w[\text{db}, \text{db}] * 2) + \\ & + w[\text{web}, \text{web}] * 3 + w[\text{db}, \text{web}] * 2 + \\ & + w[\text{web}, \text{db}] * 3 + w[\text{db}, \text{db}] * 1 + \\ & + w[\text{web}, \text{web}] * 3 + w[\text{db}, \text{web}] * 1 \end{aligned}$$

WHY?
MOTIVATING THE
(GENERAL) MODEL...

[PODZIMEK15]: COLOCATING CPU-INTENSIVE BENCHMARKS HAS SEVERE PERFORMANCE IMPACT

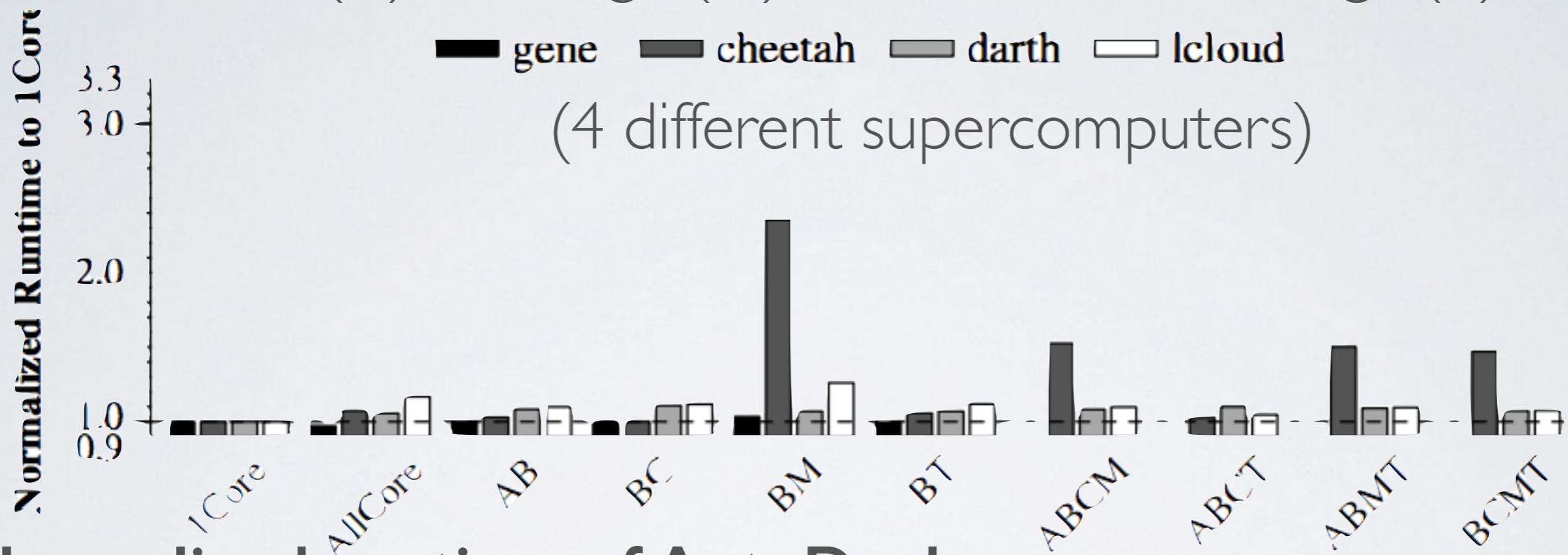
throughput of Scalac benchmark

linux kernel
allocates
workloads to
cores

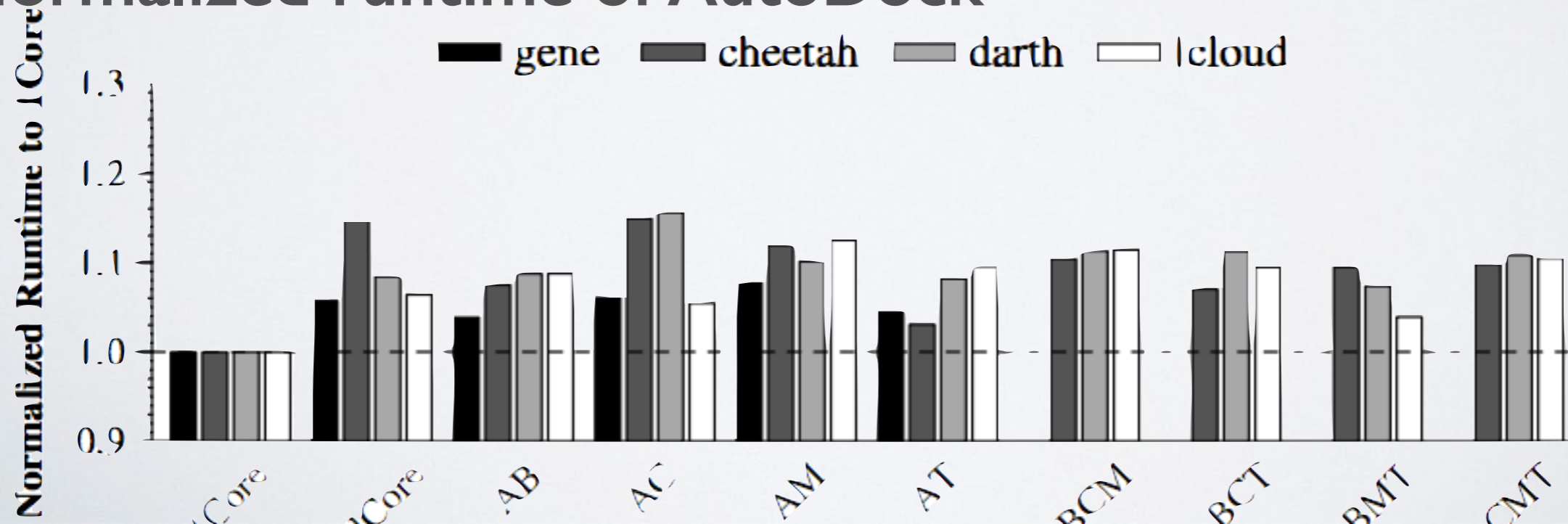


[KIM15]: PERFORMANCE IMPACT OF COLOCATION DEPENDS ON WORKLOAD

Normalized runtime of Blast when colocated with AutoDock (A), CacheBench (C), Montage (M) and ThreeKaonOmega (T)



Normalized runtime of AutoDock



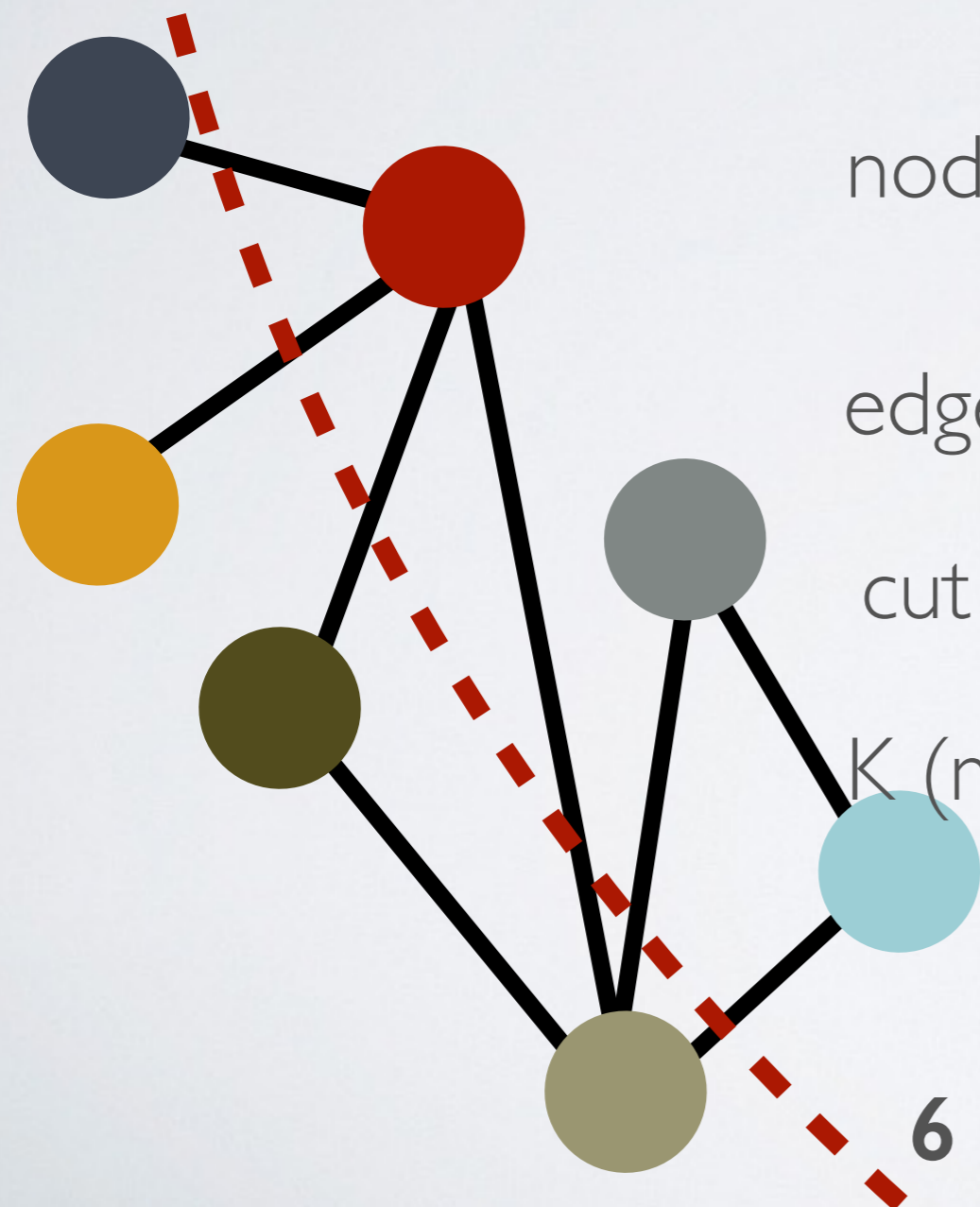
OUR RESULTS:

LINEAR COST MODEL

$$\text{COST} = \sum_{T:\text{TYPE}} \text{COEFF}_T * \text{LOAD}_T$$

THE TOTAL COST IS NP-HARD IF THERE ARE MANY TYPES

reduction from Simple Max Cut: cut a graph in two so that at least K edges cut



node — task (unit weight)

each task is of different type

edge (i,j) — $w[i,j]=1/2$; no edge — $w[k,i]=0$

cut — partition into 2 machines

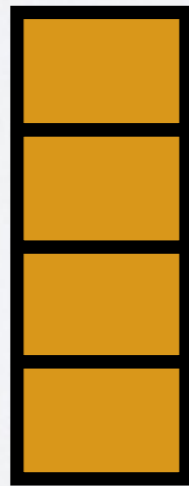
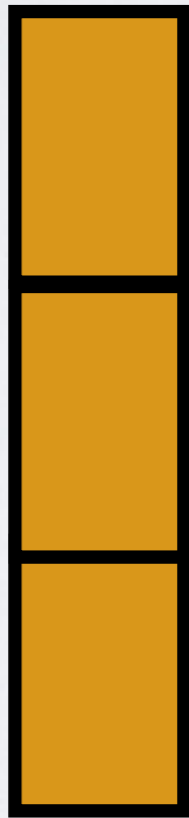
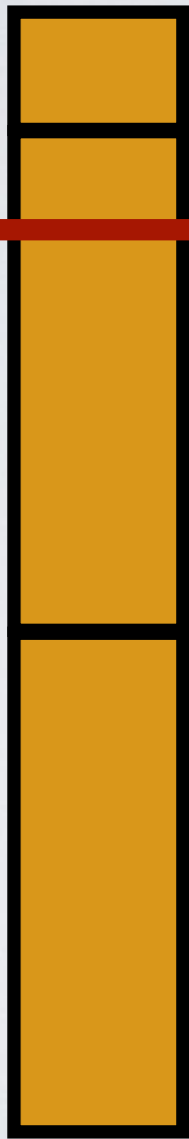
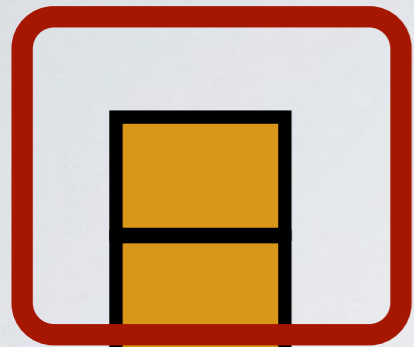
K (min # of edges cut) — max cost = $|E|-K$

6 edges cut

IN AN OPTIMAL PARTITION,
TASKS OF EACH TYPE
ARE ORDERED BY LOAD (SPT)

non-SPT assignment

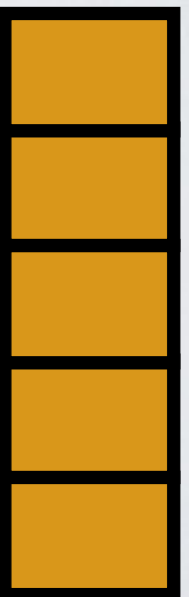
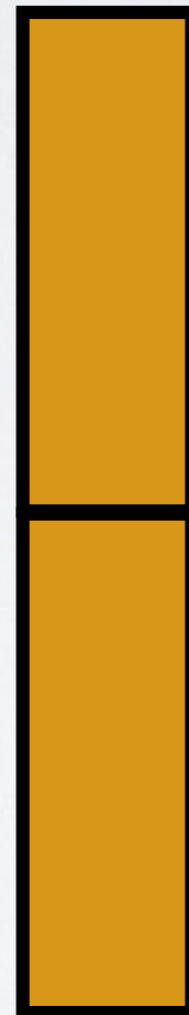
SPT assignment



M1

M2

M3



M1

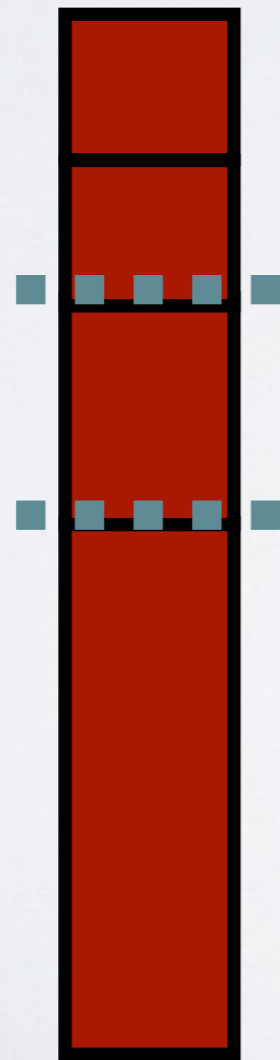
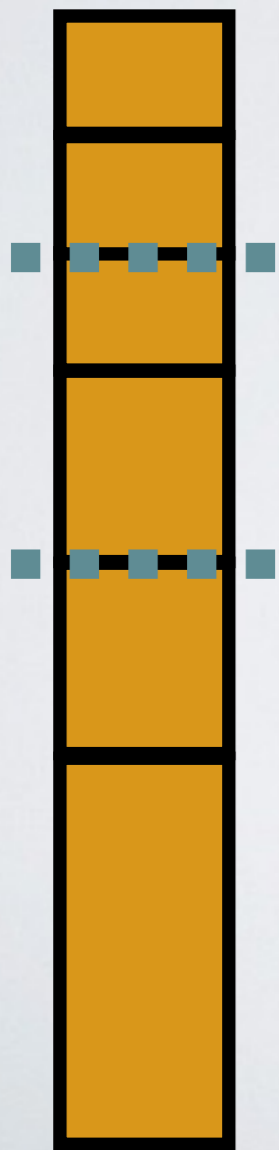
M2

M3

CUT-JUXTAPOSE: POLY IN #TASKS EXPONENTIAL IN #MACHINES AND #TYPES

1. Pick $(m-1)$ **cut** points
(each type independently)

2. **juxtapose**: find the optimal
combination (test all possibilities)



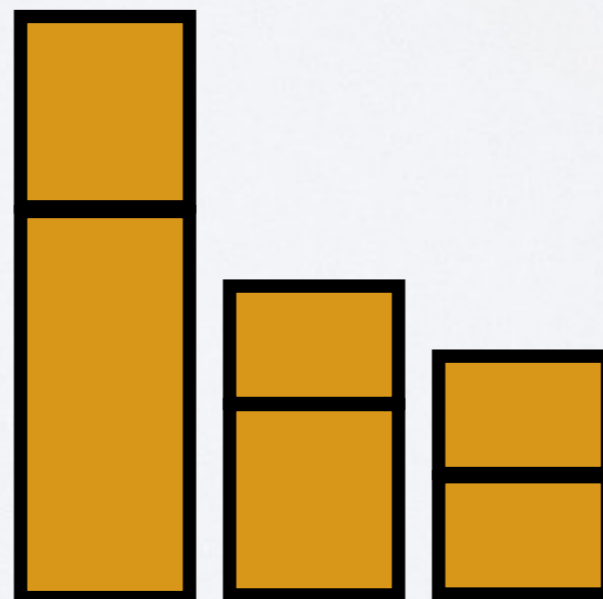
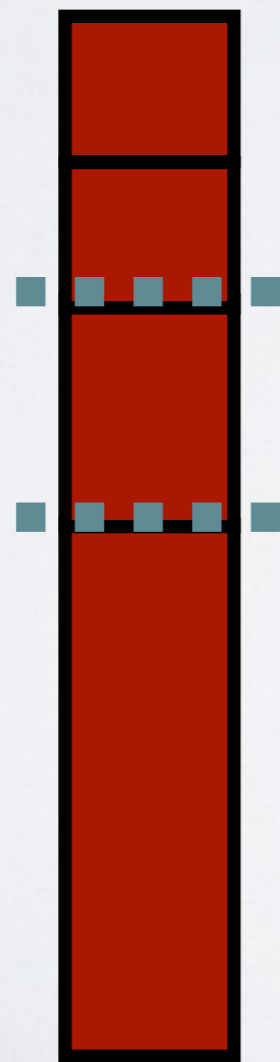
M1 M2 M3

M1 M2 M3

CUT-JUXTAPOSE: POLY IN #TASKS EXPONENTIAL IN #MACHINES AND #TYPES

1. Pick $(m-1)$ **cut** points
(each type independently)

2. **juxtapose**: find the optimal
combination (test all possibilities)



M1

M2

M3

M1

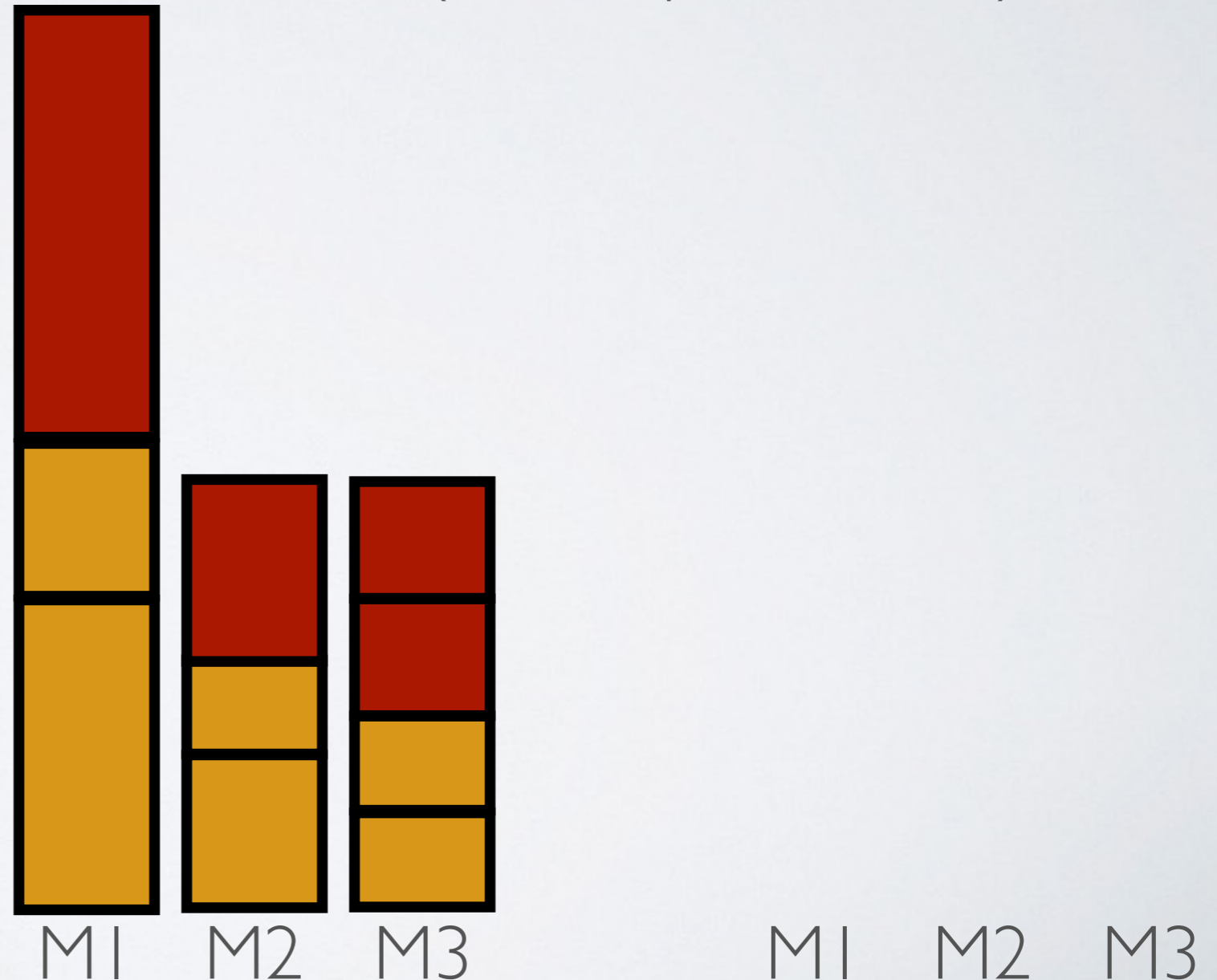
M2

M3

CUT-JUXTAPOSE: POLY IN #TASKS EXPONENTIAL IN #MACHINES AND #TYPES

1. Pick $(m-1)$ **cut** points
(each type independently)

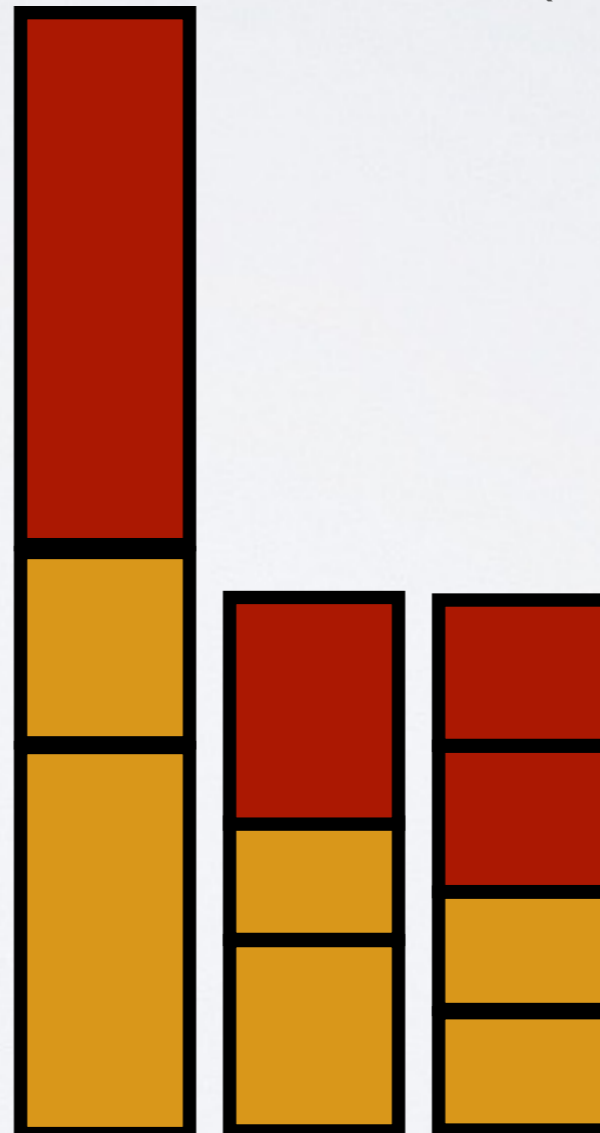
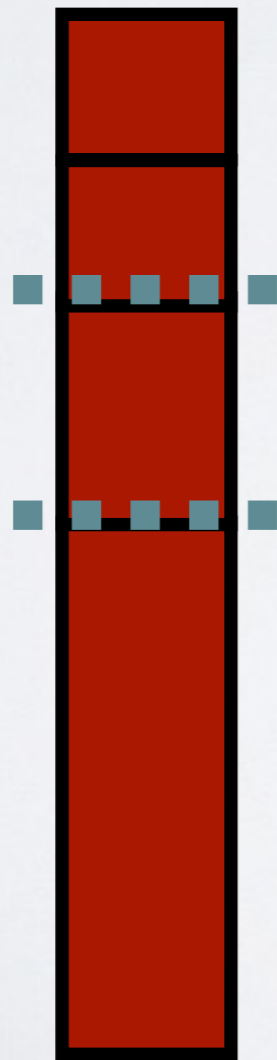
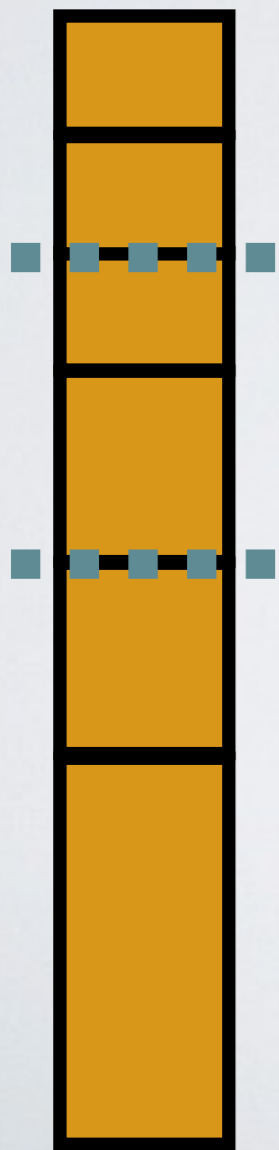
2. **juxtapose**: find the optimal
combination (test all possibilities)



CUT-JUXTAPOSE: POLY IN #TASKS EXPONENTIAL IN #MACHINES AND #TYPES

1. Pick $(m-1)$ **cut** points
(each type independently)

2. **juxtapose**: find the optimal
combination (test all possibilities)



M1

M2

M3

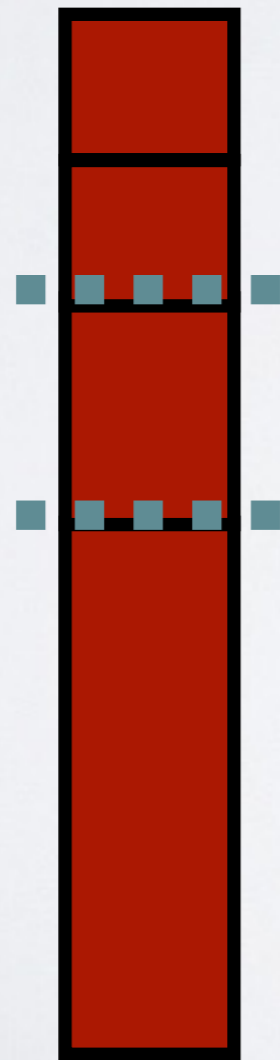
M1

M2

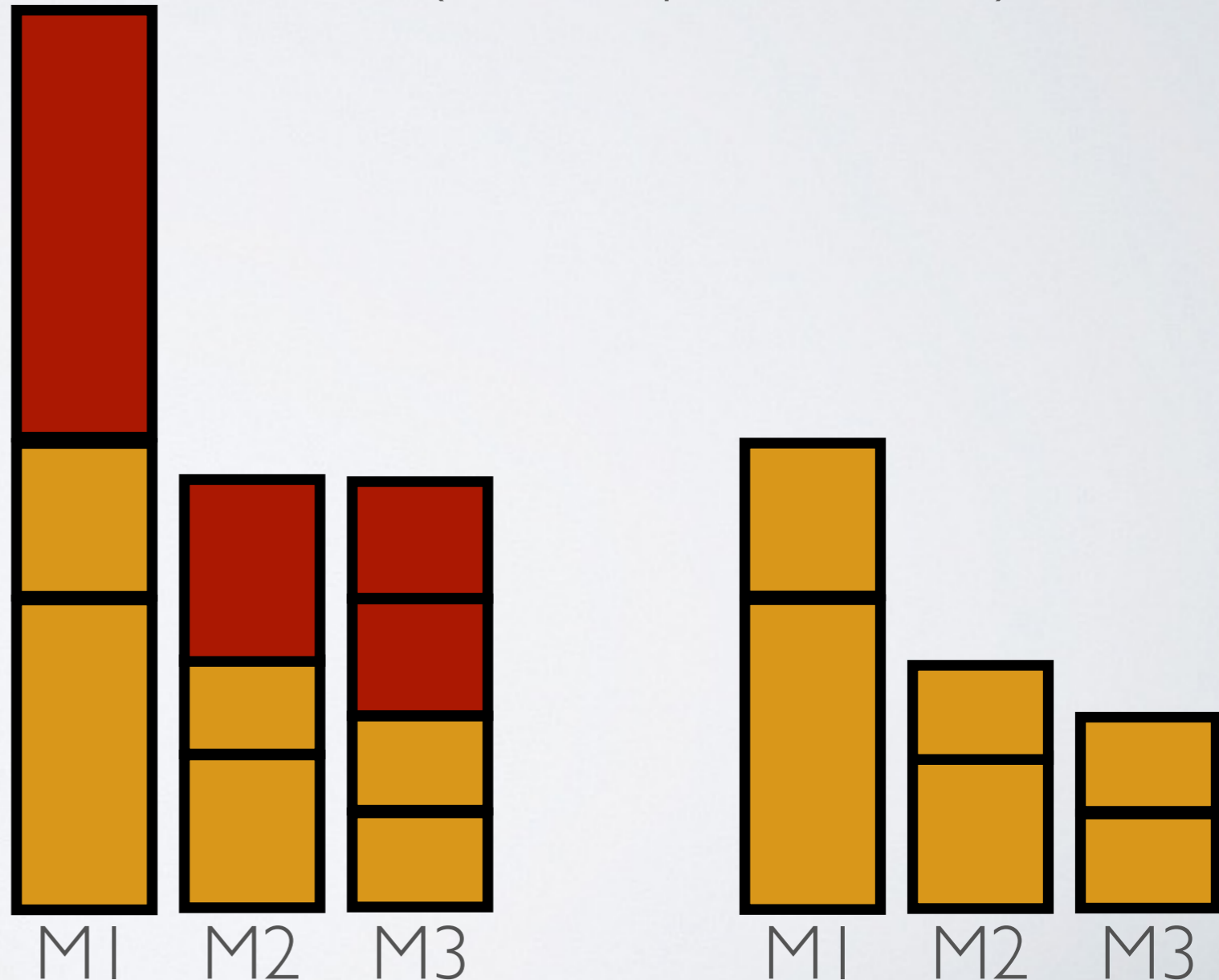
M3

CUT-JUXTAPOSE: POLY IN #TASKS EXPONENTIAL IN #MACHINES AND #TYPES

1. Pick $(m-1)$ **cut** points
(each type independently)



2. **juxtapose**: find the optimal combination (test all possibilities)



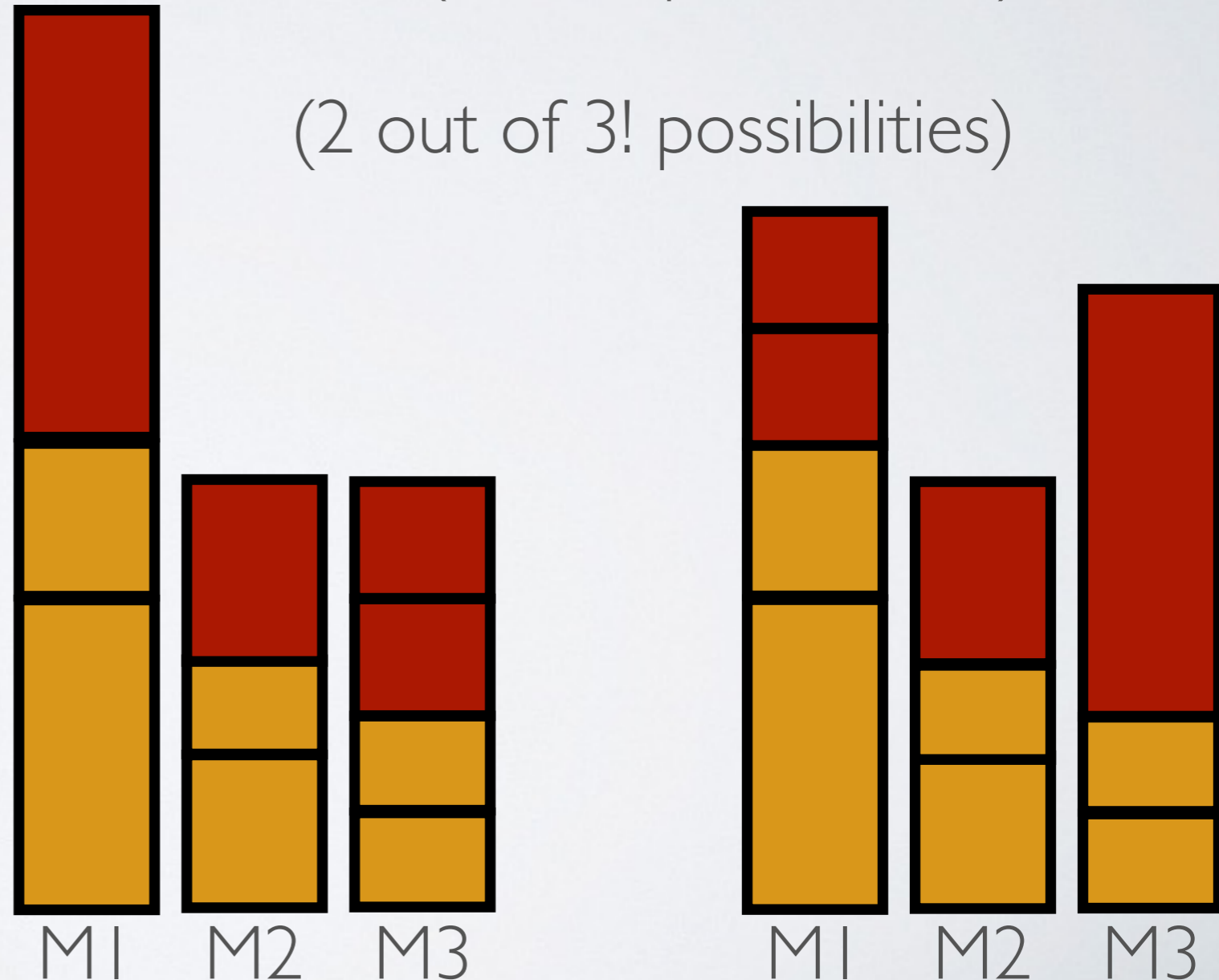
CUT-JUXTAPOSE: POLY IN #TASKS EXPONENTIAL IN #MACHINES AND #TYPES

1. Pick $(m-1)$ **cut** points
(each type independently)

2. **juxtapose**: find the optimal
combination (test all possibilities)

(2 out of 3! possibilities)

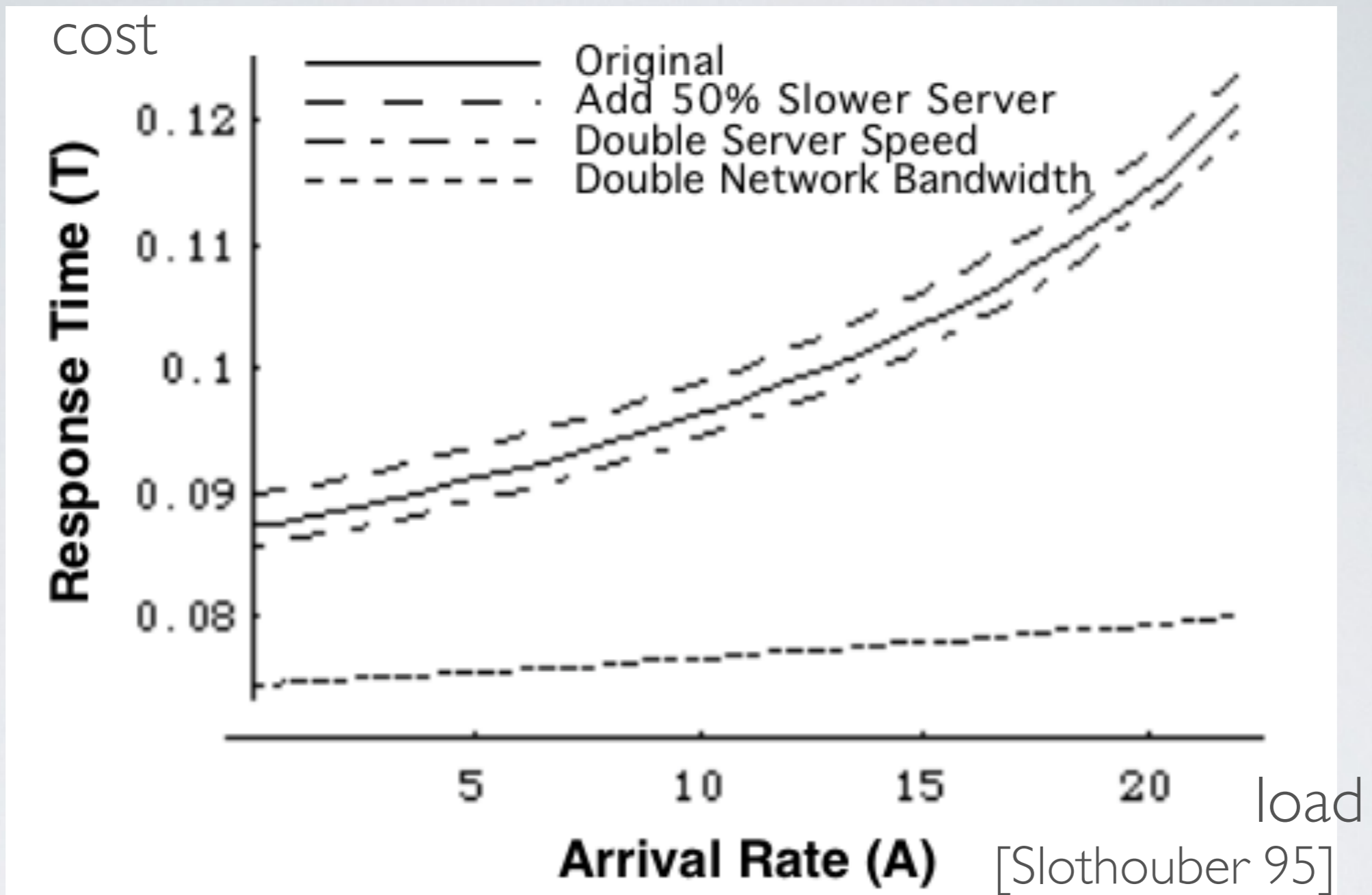
overall complexity:
 $O(n^{(m-1)T} (m!)^{T-1})$



OTHER RESULTS IN THE LINEAR COST MODEL

- dynamic programming algorithm when the number of lengths of jobs is constant $O(mn^2 \sum_t l_t)$

- dynamic programming algorithm for a single type $O(mn^2)$



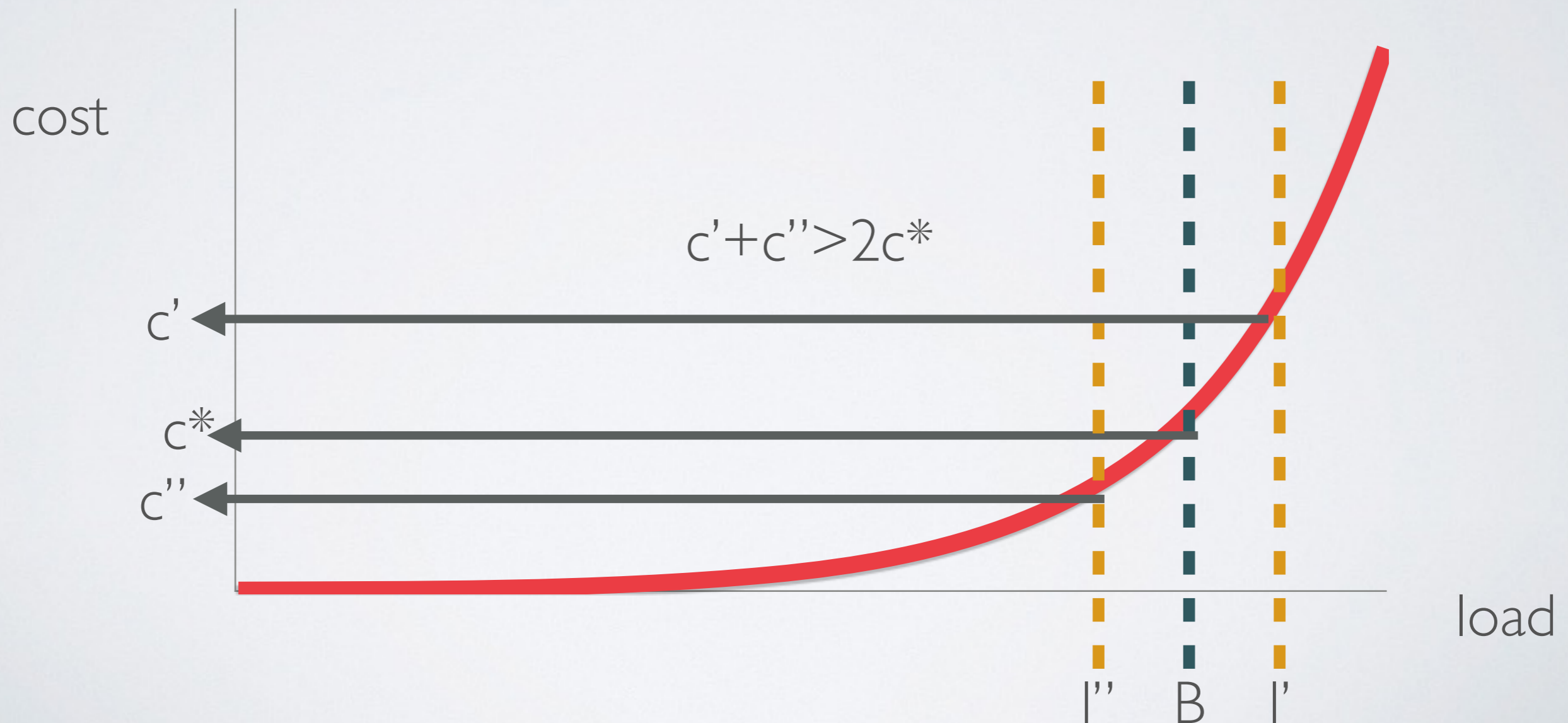
OUR RESULTS:
CONVEX COST FUNCTIONS

THE TOTAL COST IS NP-HARD IF COST FUNCTION IS STRICTLY CONVEX (EVEN FOR A SINGLE TYPE)

reduction from 3-Partition

sketch: if loads are not equal to B,

the cost of exceeding B is larger than what we save elsewhere



SUMMARY: TASKS' **TYPES** ARE USEFUL FOR ALTERNATIVE MODELING OF DATACENTER RESOURCE MANAGEMENT

- datacenters are not supercomputers! (co-allocation, more regular load, no/limited queue,)
- theoretically-sound results are rare (compared to, e.g., standard high-performance computing)
- tasks' types model tasks' heterogeneity (a webserver, a database; a computational job) and their mutual performance impact
- we have early results on complexity (few types, few machines -> poly; many types -> NP-hard); but no approx algorithms (yet?)

ACKNOWLEDGEMENTS

- Joint work with Fanny Pascual (LIP6, Paris-Sorbonne)
- Sponsored by a Google Faculty Research Award and a grant from Polish National Research Center
- Inspired by talks with Jarek Kuśmierek (Google), Piotr Skowron (University of Warsaw → Google → Oxford)
- ...and early results from services' performance measurements done by Andrzej Skrodzki (University of Warsaw)
- HiPC reviewers pointed out 2 problems in proofs!

PARTITION WITH SIDE EFFECTS

THANKS!

Krzysztof Rządca
krz@mimuw.edu.pl