

Introduction to Scheduling Theory

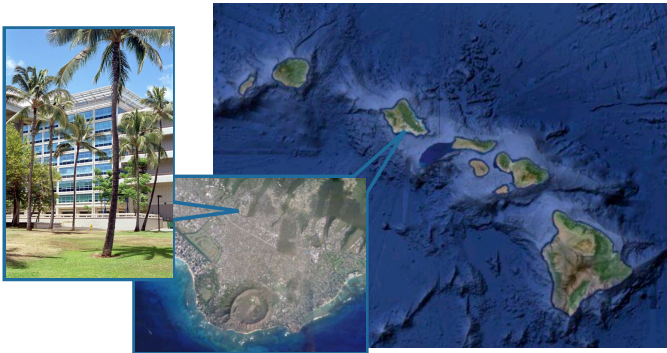
Henri Casanova^{1,2}

¹Associate Professor
Department of Information and Computer Science
University of Hawai'i at Manoa, U.S.A.

²Visiting Associate Professor
National Institute of Informatics, Japan

NII Seminar Series, October 2013

Presentation and thanks



Thanks to NII for inviting me to teach this seminar series!

Some of my research topics in last 5 years

- Scheduling (in a broad sense)
 - Divisible Load Scheduling
 - Scheduling checkpoints for fault-tolerance
 - Resource allocation in virtualized environments
 - Scheduling mixed parallel applications
 - Scheduling applications on volatile resources
 - Scheduling for energy savings
 - ...
- Simulation of distributed systems
 - Simulation tools and methodologies (SIMGRID)
 - Models for network simulation
- Random Network Topologies (with NII researchers)

Seminar topics

- Scheduling
 - A long-studied theoretical subject with practical applications
 - Comes in (too) many flavors
 - We'll explore some of them in this seminar series
- Simulation of distributed platforms and applications
 - Necessary for research on scheduling and other topics
 - Unclear and disappointing state-of-the-art
 - The SIMGRID project

Seminar organization

- Introduction to Scheduling Theory
- Scheduling Case Study: Divisible Load Scheduling
- Scheduling Case Study: Scheduling Checkpoints
- Scheduling Case Study: Scheduling Sensor Data Retrieval
- Fast and Accurate Network Simulations
- Simulating Distributed Applications with SIMGRID

Disclaimer on organization

- There are many possible topics here, especially in the area of scheduling
 - e.g., I picked 3 particular case studies but I'll likely refer to other scheduling domains as well
- I may have too much material for some topics, in which case I'll skip part of it. But my slides will of course be available to all

What is scheduling?

- Scheduling is studied in Computer Science and Operations Research
- Broad definition: *the temporal allocation of activities to resources to achieve some desirable objective*
- Examples:
 - Assign workers to machines in an factory to increase productivity
 - Pick classrooms for classes at a university to maximize the number of free classrooms on Fridays
 - Assign users to a pay-per-hour telescope to maximize profit
 - **Assign computation to processors and communications to network links so as to minimize application execution time**

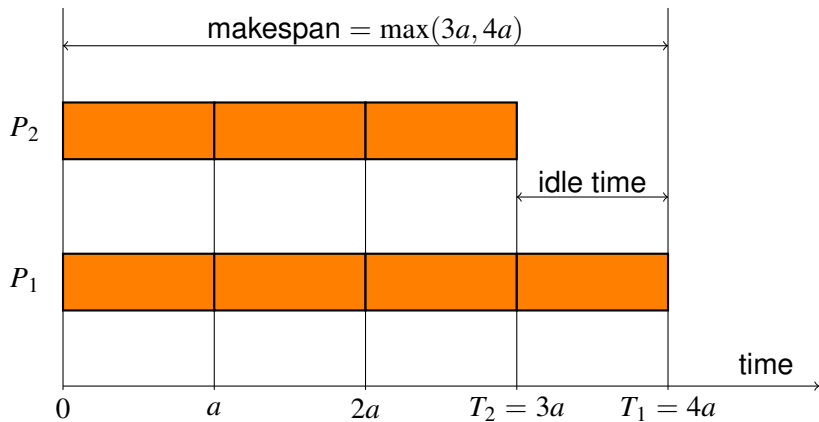
A simple scheduling problem

- A Scheduling Problem is defined by three components:
 - 1 A description of a set of resources
 - 2 A description of a set of tasks
 - 3 A description of a desired objective
- Let us get started with a simple problem: INDEP(2)
 - 1 Two identical processors, P_1 and P_2
 - Each processor can run only one task at a time
 - 2 n compute tasks
 - Each task can run on either processor in a seconds
 - Tasks are *independent*: can be computed in any order
 - 3 objective: minimize $\max(T_1, T_2)$
 - T_i is the time at which processor P_i finishes computing

The easy case

- If all tasks are *identical*, i.e., take the same amount of compute time, then the solution is obvious: Assign $\lceil n/2 \rceil$ tasks to P_1 and $\lfloor n/2 \rfloor$ tasks to P_2
 - Rule of thumb: try to have both processors finish at the same time
- The problem size is $O(1)$, the “scheduling algorithm” is $O(1)$, therefore we have a polynomial time (in fact linear) algorithm
 - For each task pick one of the two processors by comparing the index of the task with $n/2$
- We declare the problem “solved”

Gantt chart for INDEP(2) with 5 identical tasks



Non-identical tasks

- Task T_i , $i = 1, \dots, n$ takes time $a_i \geq 0$
- We say a problem is “easy” when we have a polynomial-time (p-time) algorithm:
 - Number of elementary operations is $O(f(n))$, where f is a polynomial and n is the problem size
- \mathcal{P} is the set of problems that can be solved with a p-time algorithm
- Question: is there a p-time algorithms to solve INDEP(2)?
- Disclaimer: Some of you may be familiar with algorithms and computational complexity, so bear with me while I review some fundamental background

Decision vs. optimization problem

- Complexity theory is for *decision problems*, i.e., problems that have a yes/no answer
- Scheduling problems are optimization problems
- Decision version of INDEP(2): for an integer k is there a schedule whose makespan is lower than k
- If we have a p-time algorithm for the optimization problem, then we have p-time algorithm for the decision problem
 - Run the optimization algorithm, and check whether the makespan is lower than k

Decision vs. optimization problem

- If the decision problem is in \mathcal{P} , then there is often (not always!) a p-time algorithm to solve the optimization problem
 - Binary search for the lowest k ($k \leq n \times \max_i a_i$)
 - Adds a $\log(n \times \max_i a_i)$ complexity factor, still p-time if the a_i 's are bounded (reasonable assumption)
- Almost always the case in scheduling, and decision and optimization problems are often thought of as interchangeable

Problem size?

- One has to be careful when defining the problem size
- For INDEP(2):
 - We need to enumerate n integers (the a_i 's), so the size is at least polynomial in n
 - Each a_i must be encoded (in binary) in $\lceil \log(a_i) \rceil$ bits
 - The data is $O(f(n) + \sum_{i=1}^n \lceil \log(a_i) \rceil)$, where f is a polynomial
- A problem is in P only if an algorithm exist that is polynomial in the data size as defined above

Pseudo-polynomial algorithm

- It is often possible to find algorithms polynomial in a quantity that is exponential in the (real) problem size
- For instance, to solve INDEP(2), one can resort to dynamic programming to obtain an algorithm with complexity $O(n \times \sum_{i=1}^n a_i)$
- This is a polynomial algorithm if the a_i are encoded in unary, i.e., polynomial in the numerical value of the a_i 's
- But with the a_i encoded in binary, $\sum_{i=1}^n a_i$ is exponential in the problem size!
 - To a log, linear is exponential ☺
- We say that this algorithm is *pseudopolynomial*

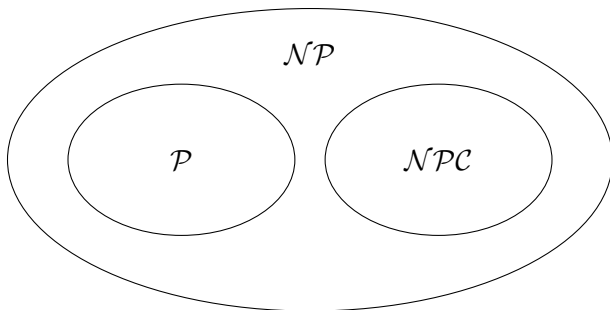
So, is INDEP(2) difficult?

- Nobody knows a p-time algorithm for solving INDEP(2)
- We define a new complexity class, \mathcal{NP}
 - Problems for which we can verify a *certificate* in p-time.
 - “Given a possible solution, can we check that the problem’s answer is Yes in p-time?”
- There are problems not in \mathcal{NP} , but not frequent
- Obviously $\mathcal{P} \subseteq \mathcal{NP}$
 - empty certificate, just solve the problem
- Big question: is $\mathcal{P} \neq \mathcal{NP}$?
 - Most people believe so, but we have no proof
 - For all the follows, “unless $\mathcal{P} = \mathcal{NP}$ ” is implied

\mathcal{NP} -complete problems

- Some problems in \mathcal{NP} are at least as difficult as all other problems in \mathcal{NP}
- They are called \mathcal{NP} -complete, and their set is \mathcal{NPC}
- Cook's theorem: The SAT problems is in \mathcal{NPC}
 - Satisfiability of a boolean conjunction of disjunctions
- How to prove that a problem, P , is \mathcal{NP} -complete:
 - Prove that $P \in \mathcal{NP}$ (typically easy)
 - Prove that P *reduces* to Q , where $Q \in \mathcal{NPC}$ (can be hard)
 - For an instance I_Q construct in p-time an instance I_P
 - Prove that I_P has a solution if and only if I_Q has a solution
- By now we know many problems in \mathcal{NPC}
- Goal: pick $Q \in \mathcal{NPC}$ so that the reduction is easy

Well-known complexity classes



INDEP(2) is \mathcal{NP} -complete

- INDEP(2) (decision version) is in \mathcal{NP}
 - Certificate: for each a_i whether it is schedule on P_1 or P_2
 - In linear time, compute the makespan on both processors, and compare to k to answer "Yes"
- Let us consider an instance of 2-PARTITION $\in \mathcal{NPC}$:
 - Given n integers x_i , is there a subset I of $\{1, \dots, n\}$ such that $\sum_{i \in I} x_i = \sum_{i \notin I} x_i$?
- Let us construct an instance of INDEP(2):
 - Let $k = \frac{1}{2} \sum x_i$, let $a_i = x_i$
- The proof is trivial
 - If k is non-integer, neither instance has a solution
 - Otherwise, each processor corresponds to one subset
- In fact, INDEP(2) is essentially identical to 2-PARTITION

So what?

- This \mathcal{NP} -completeness proof is probably the most trivial in the world 😊
- But now we are thus pretty sure that there is no p-time algorithm to solve INDEP(2)
- What we look for now are *approximation algorithms*...

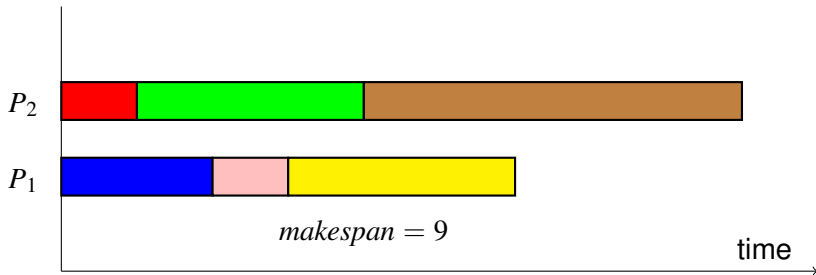
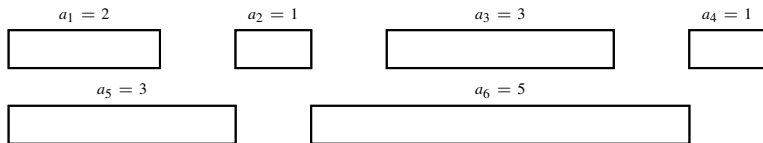
Approximation algorithms

- Consider an optimization problem
- A p-time algorithm is a λ -*approximation algorithm* if it returns a solution that's at most a factor λ from the optimal solution (the closer λ to 1, the better)
 - λ is called the *approximation ratio*
- *Polynomial Time Approximation Scheme* (PTAS): for any ϵ there exists a $(1 + \epsilon)$ -approximation algorithm (may be non-polynomial is $1/\epsilon$)
- *Fully Polynomial Time Approximation Scheme* (FPTAS): for any ϵ there exists a $(1 + \epsilon)$ -approximation algorithm polynomial in $1/\epsilon$
- Typical goal: find a FPTAS, if not find a PTAS, if not find a λ -approximation for a low value of λ

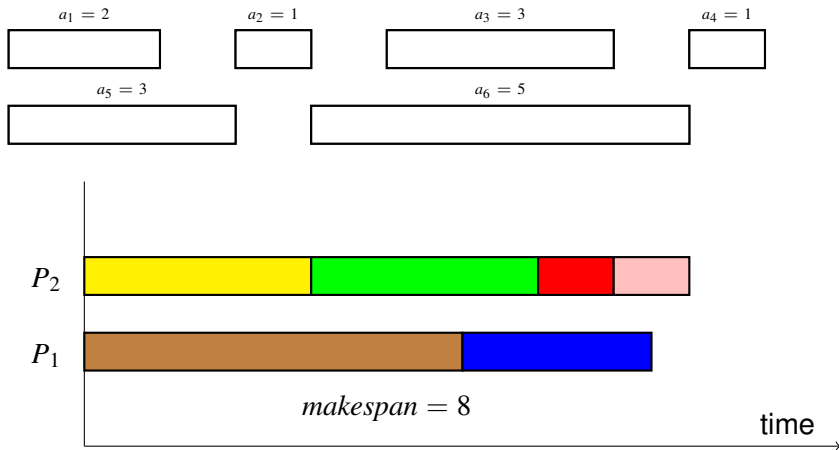
Greedy algorithms

- A greedy algorithm is one that builds a solution step-by-step, via local incremental decisions
- It turns out that several greedy scheduling algorithms are approximation algorithms
 - Informally, they're not as "bad" as one may think
- Two natural greedy algorithms for INDEP(2):
 - **greedy-online**: take the tasks in arbitrary order and assign each task to the least loaded processor
 - We don't know which tasks are coming
 - **greedy-offline**: sort the tasks by decreasing a_i , and assign each task in that order to the least loaded processor
 - We know all the tasks ahead of time

Example with 6 tasks: Online



Example with 6 tasks: Offline



Greedy-online for INDEP(2)

Theorem

Greedy-online is a $\frac{3}{2}$ -approximation

■ Proof:

- P_i finishes computing at time M_i (M stands for makespan)
- Let us assume $M_1 \geq M_2$ ($M_{\text{greedy}} = M_1$)
- Let T_j the last task to execute on P_1
- Since the greedy algorithm put T_j on P_1 , then $M_1 - a_j \leq M_2$
- We have $M_1 + M_2 = \sum_i a_i = S$
- $M_{\text{greedy}} = M_1 = \frac{1}{2}(M_1 + (M_1 - a_j) + a_j) \leq \frac{1}{2}(M_1 + M_2 + a_j) = \frac{1}{2}(S + a_j)$
- but $M_{\text{opt}} \geq S/2$ (ideal lower bound on optimal)
- and $M_{\text{opt}} \geq a_j$ (at least one task is executed)
- Therefore: $M_{\text{greedy}} \leq \frac{1}{2}(2M_{\text{opt}} + M_{\text{opt}}) = \frac{3}{2}M_{\text{opt}} \quad \square$

Greedy-offline for INDEP(2)

Theorem

Greedy-offline is a $\frac{7}{6}$ -approximation

■ Proof:

- If $a_j \leq \frac{1}{3}M_{opt}$, the previous proof can be used
 - $M_{greedy} \leq \frac{1}{2}(2M_{opt} + \frac{1}{3}M_{opt}) = \frac{7}{6}M_{opt}$
- If $a_j > \frac{1}{3}M_{opt}$, then $j \leq 4$
 - if T_j was the 5th task, then, due to the task ordering, there would be 5 tasks with $a_i > \frac{1}{3}M_{opt}$
 - There would be at least 3 tasks on the same processor in the optimal schedule
 - Therefore $M_{opt} > 3 \times \frac{1}{3}M_{opt}$, a contradiction
- One can check all possible scenarios for 4 tasks and show optimality □

Bounds are tight

■ Greedy-online:

- a_i 's = $\{1, 1, 2\}$
- $M_{greedy} = 3; M_{opt} = 2$
- $ratio = \frac{3}{2}$

■ Greedy-offline:

- a_i 's = $\{3, 3, 2, 2, 2\}$
- $M_{greedy} = 7; M_{opt} = 6$
- $ratio = \frac{7}{6}$

PTAS and FPTAS for INDEP(2)

Theorem

There is a PTAS $((1 + \epsilon)$ -approximation) for INDEP(2)

■ Proof Sketch:

- Classify tasks as either “small” or “large”
 - Very common technique
- Replace all small tasks by same-size tasks
- Compute an optimal schedule of the modified problem in p-time (not polynomial in $1/\epsilon$)
- Show that the cost is $\leq 1 + \epsilon$ away from the optimal cost
- The proof is a couple of pages, but not terribly difficult

Theorem

There is a FPTAS $((1 + \epsilon)$ -approx pol. in $1/\epsilon$) for INDEP(2)

We know a lot about INDEP(2)

- INDEP(2) is NP-complete
 - We have simple greedy algorithms with guarantees on result quality
 - We have a simple PTAS
 - We even have a (less simple) FPTAS
 - INDEP(2) is basically "solved"
-
- Sadly, not many scheduling problems are this well-understood...

INDEP(P) is much harder

- INDEP(P) is \mathcal{NP} -complete by trivial reduction to 3-PARTITION:
 - Give $3n$ integers a_1, \dots, a_{3n} and an integer B , can we partition the $3n$ integers into n sets, each of sum B ? (assuming that $\sum_i a_i = nB$)
- 3-PARTITION is \mathcal{NP} -complete “in the strong sense”, unlike 2-PARTITION
 - Even when encoding the input in unary (i.e., no logarithmic numbers of bits), one cannot find an algorithm polynomial in the size of the input!
 - Informally, a problem is \mathcal{NP} -complete “in the weak sense” if it is hard only if the numbers in the input are unbounded
- INDEP(P) is thus fundamentally harder than INDEP(2)

Approximation algorithm for INDEP(p)

Theorem

Greedy-online is a $(2 - \frac{1}{p})$ -approximation

■ Proof (usual reasoning):

- Let $M_{\text{greedy}} = \max_{1 \leq i \leq p} M_i$, and j be such that $M_j = M_{\text{greedy}}$
- Let T_k be the last task assigned to processor P_j
- $\forall i, \quad M_i \geq M_j - a_k$ (greedy algorithm)
- $S = \sum_i^p M_i = M_j + \sum_{i \neq j} M_i \geq M_j + (p-1)(M_j - a_k) = pM_j + (p-1)a_k$
- Therefore, $M_{\text{greedy}} = M_j \leq \frac{S}{p} + (1 - \frac{1}{p})a_k$
- But $M_{\text{opt}} \geq a_k$ and $M_{\text{opt}} \geq S/p$
- So $M_{\text{greedy}} \leq M_{\text{opt}} + (1 - \frac{1}{p})M_{\text{opt}} \quad \square$
- This ratio is “tight” (e.g., an instance with $p(p-1)$ tasks of size 1 and one task of size p has this ratio)

Approximation algorithm for INDEP(p)

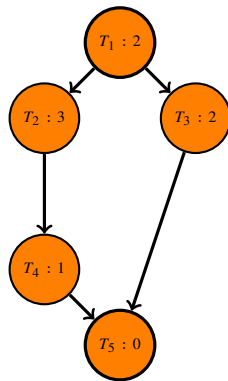
Theorem

Greedy-offline is a $(\frac{4}{3} - \frac{1}{3p})$ -approximation

- The proof is more involved, but follows the spirit of the proof for INDEP(2)
- This ratio is tight
- There is a PTAS for INDEP(p), a $(1 + \epsilon)$ -approximation (massively exponential in $1/\epsilon$)
- There is no known FPTAS, unlike for INDEP(2)

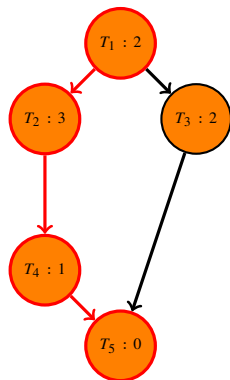
Task dependencies

- In practice tasks often have *dependencies*
- A general model of computation is the Acyclic Directed Graph (DAG), $G = (V, E)$
- Each task has a *weight* (i.e., execution time in seconds), a *parent*, and *children*
- The first task is the *source*, the last task the *sink*
- Topological (partial) order of the tasks



Critical path

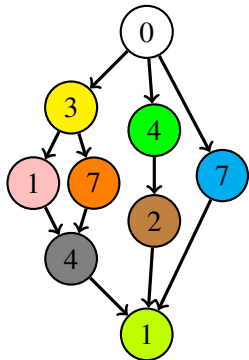
- Assume that the DAG executes on p processors
- The longest path (in seconds) is called the *critical path*
- The length of the critical path (CP) is a *lower bound on M_{opt}* , regardless of the number of processors
- In this example, the CP length is 6 (the other path has length 4)



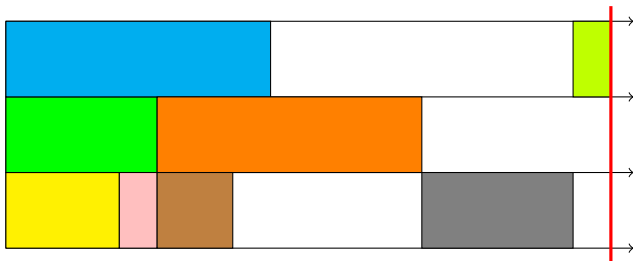
Complexity

- Unsurprisingly, DAG scheduling is \mathcal{NP} -complete
 - Independent tasks is a special case of DAG scheduling
- Typical greedy algorithm skeleton:
 - Maintain a list of *ready* tasks (with cleared dependencies)
 - Greedily assign a ready task to an available processor as early as possible (don't leave a processor idle unnecessarily)
 - Update the list of ready tasks
 - Repeat until all tasks have been scheduled
- This is called **List Scheduling**
- Many list scheduling algorithms are possible
 - Depending on how to select the ready task to schedule next

List scheduling example



3 Processors



Makespan = 16; CP Length = 15

Idle Time = $1+5+5+8 = 19$

List scheduling

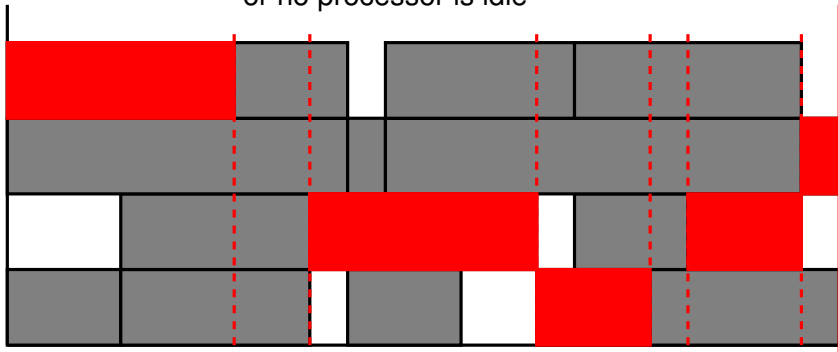
Theorem (fundamental)

List scheduling is a $(2 - \frac{1}{p})$ -approximation

- Doesn't matter how the next ready task is selected
- Let's prove this theorem informally
 - Really simple proof if one doesn't use the typical notations for schedules
 - I never use these notations in public 😊

Approximation ratio

At any point in time either a task on the red path is running or no processor is idle



Approximation ratio

- Let L be the length of the red path (in seconds), p the number of processors, I the total idle time, M the makespan, and S the sum of all task weights
 - $I \leq (p - 1)L$
 - processors can be idle only when a red task is running
 - $L \leq M_{opt}$
 - The optimal makespan is longer than any path in the DAG
 - $M_{opt} \geq S/p$
 - S/p is the makespan with zero idle time
 - $p \times M = I + S$
 - rectangle's area = white boxes + non-white boxes
- $\Rightarrow p \times M \leq (p - 1)M_{opt} + pM_{opt} \Rightarrow M \leq (2 - \frac{1}{p})M_{opt} \quad \square$

Good list scheduling?

- All list scheduling algorithms thus have the same approximation ratio
- But there are many options for list scheduling
 - Many ways of sorting the ready tasks...
- In practice, some may be better than others
- One well-known option, *Critical path scheduling*

Critical path scheduling

- When given a set of ready tasks, which one do we pick to schedule?
- Idea: pick a task on the CP
 - If we prioritize tasks on the CP, then the CP length is reduced
 - The CP length is a lower bound on the makespan
 - So intuitively it's good for it to be low
- For each (ready) task, compute its *bottom level*, the length of the path from the task to the sink
- Pick the task with the *largest* bottom level

Graham's notation

- There are SO many variations on the scheduling problem that Graham has proposed a standard notation: $\alpha|\beta|\gamma$
 - *alpha*: processors
 - *beta*: tasks
 - *gamma*: objective function
- Let's see some examples for each

α : processors

- 1: one processor
- Pn : n identical processors (if n not fixed, not given)
- Qn : n uniform processors (if n not fixed, not given)
 - Each processor has a (different) compute speed
- Rn : n unrelated processors (if n not fixed, not given)
 - Each processor has a (different) compute speed for each (different) task (e.g., P_1 can be faster than P_2 for T_1 , but slower for T_2)

β : tasks

- r_j : tasks have *release dates*
- d_j : tasks have *deadlines*
- $p_j = x$: all tasks have weight x
- $prec$: general precedence constraints (DAG)
- $tree$: tree precedence constraints
- $chains$: chains precedence constraints (multiple independent paths)
- $pmtn$: tasks can be preempted and restarted (on other processors)
 - Makes scheduling easier, and can often be done in practice
- ...

γ : objective function

- C_{max} : makespan
- $\sum C_i$: mean flow-time (completion time minus release date if any)
- $\sum w_i C_i$: average weighted flow-time
- L_{max} : maximum lateness ($\max(0, C_i - d_i)$)
- ...

Example scheduling problems

- The classification is not perfect and variations among authors are common
- Some examples:
 - $P2||C_{max}$, which we called INDEP(2)
 - $P||C_{max}$, which we called INDEP(P)
 - $P|prec|C_{max}$, which we called DAG scheduling
 - $R2|chains|\sum C_i$
 - Two related processors, chains, minimize sum-flow
 - $P|r_j; p_j \in \{1, 2\}; d_j; pmtn|L_{max}$
 - Identical processors, tasks with release dates and deadlines, task weights either 1 or 2, preemption, minimize maximum lateness

Where to find known results

- Luckily, the body of knowledge is well-documented (and Graham's notation widely used)
- Several books on scheduling that list known results
 - *Handbook of Scheduling*, Leung and Anderson
 - *Scheduling Algorithms*, Brucker
 - *Scheduling: Theory, Algorithms, and Systems*, Pinedo
 - ...
- Many published survey articles

Example list of known results

■ Excerpt from
*Scheduling
Algorithm*, P.
Brucker

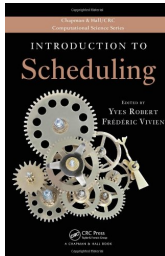
$P2 \parallel C_{max}$	Lenstra et al. [155]
* $P \parallel C_{max}$	Garey & Johnson [98]
* $P \mid p_i = 1;intree;r_i \mid C_{max}$	Brucker et al. [35]
* $P \mid p_i = 1;prec \mid C_{max}$	Ullman [203]
* $P2 \mid chains \mid C_{max}$	Du et al. [86]
* $Q \mid p_i = 1;chains \mid C_{max}$	Kubiak [129]
* $P \mid p_i = 1;outtree \mid L_{max}$	Brucker et al. [35]
* $P \mid p_i = 1;intree;r_i \mid \sum C_i$	Lenstra [150]
* $P \mid p_i = 1;prec \mid \sum C_i$	Lenstra & Rinnooy Kan [152]
* $P2 \mid chains \mid \sum C_i$	Du et al. [86]
* $P2 \mid r_i \mid \sum C_i$	Single-machine problem
$P2 \parallel \sum w_i C_i$	Bruno et al. [58]
* $P \parallel \sum w_i C_i$	Lenstra [150]
* $P2 \mid p_i = 1;chains \mid \sum w_i C_i$	Timkovsky [201]
* $P2 \mid p_i = 1;chains \mid \sum U_i$	Single-machine problem
* $P2 \mid p_i = 1;chains \mid \sum T_i$	Single-machine problem

Table 5.3: \mathcal{NP} -hard parallel machine problems without preemption.

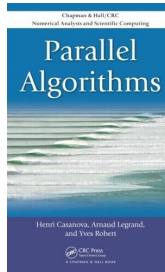
Conclusion

- Scheduling problems are diverse and often difficult
- Relevant theoretical questions:
 - Is it in \mathcal{P} ?
 - Is it \mathcal{NP} -complete?
 - Are there approximation algorithms?
 - Are there PTAS or FTPAS?
 - Are there at least decent non-guaranteed heuristics?
- Luckily, scheduling problems have been studied a lot
- Come up with the Graham notation for your problem and check what is known about it!

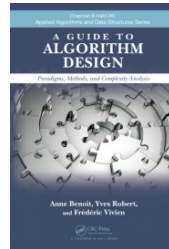
Sources and acknowledgments



Y. Robert
F. Vivien



H. Casanova
A. Legrand
Y. Robert



A. Benoit
Y. Robert
F. Vivien