

Introduction to Scheduling

Arnaud Legrand, CNRS, University of Grenoble

LIG laboratory, arnaud.legrand@imag.fr

November 9, 2009

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Matrix Product: Sequential Version

```
1 { To compute  $C \leftarrow C + A \times B$  }
2 for  $i = 1$  to  $n$  do
3   for  $j = 1$  to  $n$  do
4     for  $k = 1$  to  $n$  do
5        $C_{i,j} \leftarrow C_{i,j} + A_{i,k} \times B_{k,j}$ 
```

$B_{1,1}$	$B_{1,2}$
$B_{2,1}$	$B_{2,2}$

$A_{1,1}$	$A_{1,2}$
$A_{2,1}$	$A_{2,2}$

$C_{1,1}$	$C_{1,2}$
$C_{2,1}$	$C_{2,2}$

Matrix Product: Sequential Version

$$2 \quad C_{1,1} \leftarrow C_{1,1} + A_{1,1} \times B_{1,1}$$

$$4 \quad C_{1,1} \leftarrow C_{1,1} + A_{1,2} \times B_{2,1}$$

$$6 \quad C_{1,2} \leftarrow C_{1,2} + A_{1,1} \times B_{1,2}$$

$$8 \quad C_{1,2} \leftarrow C_{1,2} + A_{1,2} \times B_{2,2}$$

9 ...

$B_{1,1}$	$B_{1,2}$
$B_{2,1}$	$B_{2,2}$

CPU

2	4	6	8
---	---	---	---

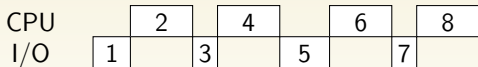
$A_{1,1}$	$A_{1,2}$
$A_{2,1}$	$A_{2,2}$

$C_{1,1}$	$C_{1,2}$
$C_{2,1}$	$C_{2,2}$

Matrix Product: Sequential Version

- 1 Load $C_{1,1}$, $A_{1,1}$, $B_{1,1}$
- 2 $C_{1,1} \leftarrow C_{1,1} + A_{1,1} \times B_{1,1}$
- 3 Unload $A_{1,1}$, $B_{1,1}$. Load $A_{1,2}$, $B_{2,1}$
- 4 $C_{1,1} \leftarrow C_{1,1} + A_{1,2} \times B_{2,1}$
- 5 Unload $C_{1,1}$, $A_{1,2}$, $B_{2,1}$. Load $C_{1,2}$, $A_{1,1}$, $B_{1,2}$
- 6 $C_{1,2} \leftarrow C_{1,2} + A_{1,1} \times B_{1,2}$
- 7 Unload $A_{1,1}$, $B_{1,2}$
- 8 $C_{1,2} \leftarrow C_{1,2} + A_{1,2} \times B_{2,2}$
- 9 ...

$B_{1,1}$	$B_{1,2}$
$B_{2,1}$	$B_{2,2}$



$A_{1,1}$	$A_{1,2}$
$A_{2,1}$	$A_{2,2}$

$C_{1,1}$	$C_{1,2}$
$C_{2,1}$	$C_{2,2}$

Matrix Product: Sequential Version

- 1 Load $C_{1,1}$, $A_{1,1}$, $B_{1,1}$
- 2 $C_{1,1} \leftarrow C_{1,1} + A_{1,1} \times B_{1,1}$
- 3 Unload $A_{1,1}$, $B_{1,1}$. Load $A_{1,2}$, $B_{2,1}$
- 4 $C_{1,1} \leftarrow C_{1,1} + A_{1,2} \times B_{2,1}$
- 5 Unload $C_{1,1}$, $A_{1,2}$, $B_{2,1}$. Load $C_{1,2}$, $A_{1,1}$, $B_{1,2}$
- 6 $C_{1,2} \leftarrow C_{1,2} + A_{1,1} \times B_{1,2}$
- 7 Unload $A_{1,1}$, $B_{1,2}$
- 8 $C_{1,2} \leftarrow C_{1,2} + A_{1,2} \times B_{2,2}$
- 9 ...

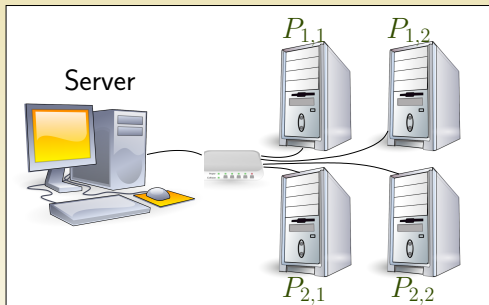
$B_{1,1}$	$B_{1,2}$
$B_{2,1}$	$B_{2,2}$

Sequential Programs

Sequential programs are generally a succession of CPU burst and I/O burst.

Matrix Product: Parallel Version (1/2)

Setting

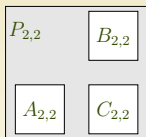
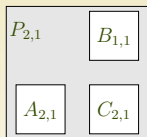
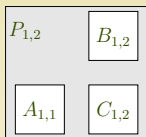
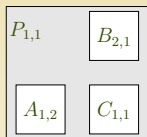


- ▶ A , B , and C are initially located on the server.
- ▶ We will distribute A , B , and C on $P_{1,1}$, $P_{1,2}$, $P_{2,1}$, $P_{2,2}$.
- ▶ We will make use of all four processors to compute $C \leftarrow C + A \times B$.
- ▶ Such a parallel program could be written using for example MPI. We want a SPMD algorithm.

Matrix Product: Parallel Version (2/2)

Algorithm

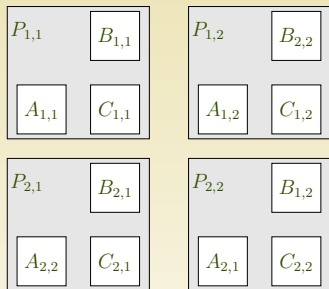
- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 Unload $C_{i,j}$ to the server



Matrix Product: Parallel Version (2/2)

Algorithm

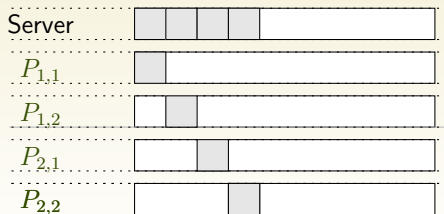
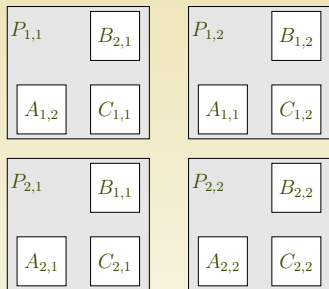
- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 Unload $C_{i,j}$ to the server



Matrix Product: Parallel Version (2/2)

Algorithm

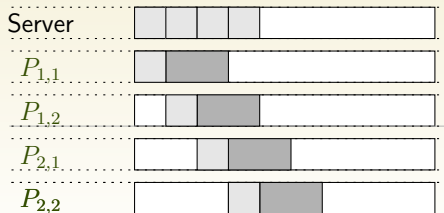
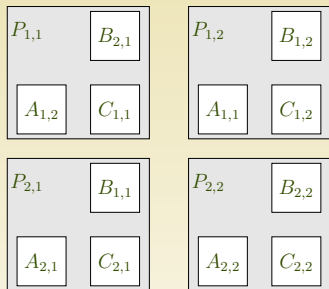
- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 Unload $C_{i,j}$ to the server



Matrix Product: Parallel Version (2/2)

Algorithm

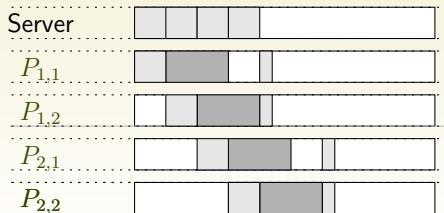
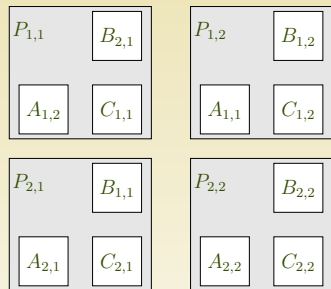
- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 Unload $C_{i,j}$ to the server



Matrix Product: Parallel Version (2/2)

Algorithm

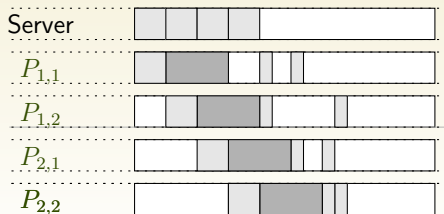
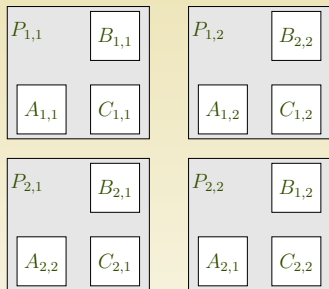
- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 Unload $C_{i,j}$ to the server



Matrix Product: Parallel Version (2/2)

Algorithm

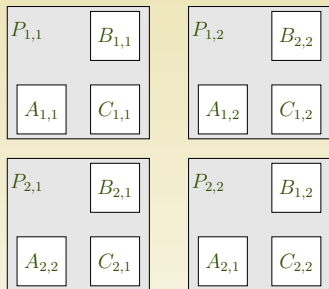
- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 Unload $C_{i,j}$ to the server



Matrix Product: Parallel Version (2/2)

Algorithm

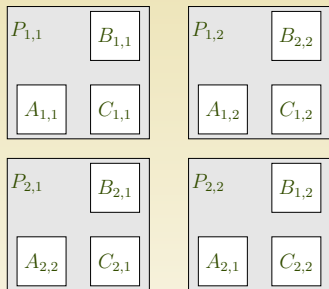
- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 Unload $C_{i,j}$ to the server



Matrix Product: Parallel Version (2/2)

Algorithm

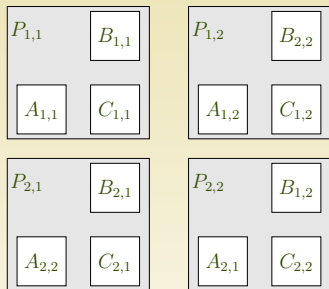
- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 **Unload $C_{i,j}$ to the server**



Matrix Product: Parallel Version (2/2)

Algorithm

- 1 $\{ P_{i,j} \text{ is responsible for computing } C_{i,j}. \}$
- 2 Load $C_{i,j}$, $A_{i,(i+j)\%2}$, $B_{(i+j)\%2,j}$ from the server
- 3 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 4 Exchange A_{local} with horizontal neighbor
- 5 Exchange B_{local} with vertical neighbor
- 6 $C_{local} \leftarrow C_{local} + A_{local} \times B_{local}$
- 7 Unload $C_{i,j}$ to the server



Server

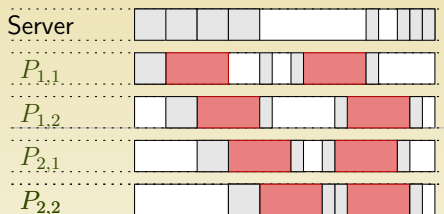
Parallel Programs

Parallel programs are generally a succession of **CPU burst** and **communication burst**. The synchronization pattern generally incurs **idle time**. This is the **parallelization overhead**.

$P_{2,2}$

Work, Cost, Speed-up and Efficiency



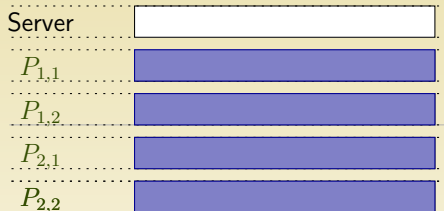


Definition: **Work**.

The **work** is the **amount of computation** performed (the surface of the pink area).

In the previous parallel Matrix Multiplication example, the work is the same as in the sequential Matrix Multiplication example.

However, parallel algorithms generally do not do the same operations as the sequential ones. They often have to do more. Therefore, the work $W(p)$ generally **depends on the number of processors that are allotted!**



Definition: **Cost**.

$$C(p) = p \times \text{TotalTime}(p).$$

It is the total surface.

The cost accounts for the idle time of the processing units.



Definition: Speed-up and Efficiency.

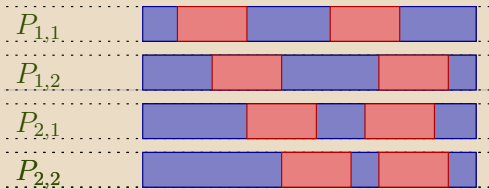
► **Speed-up:** $s(p) = \frac{\text{SequentialTime}}{\text{TotalTime}(p)}$.

► **Efficiency:** $e(p) = \frac{s(p)}{p} = \frac{\text{SequentialTime}}{p \times \text{TotalTime}(p)}$.

Side Note on Speed-up and Efficiency

Speed-up

We have $\text{SequentialTime} \leq C(p) \leq p \times \text{TotalTime}(p)$.



Speed-up

We have $\text{SequentialTime} \leq C(p) \leq p \times \text{TotalTime}(p)$.

Hence, $s(p) = \frac{\text{SequentialTime}}{\text{TotalTime}(p)} \leq p$ and $e(p) = \frac{\text{SequentialTime}}{p \text{TotalTime}(p)} \leq 1$.

The speed-up is bounded by the number of processors and the efficiency is thus in $[0, 1]$.

Speed-up

We have $\text{SequentialTime} \leq C(p) \leq p \times \text{TotalTime}(p)$.

Hence, $s(p) = \frac{\text{SequentialTime}}{\text{TotalTime}(p)} \leq p$ and $e(p) = \frac{\text{SequentialTime}}{p \text{TotalTime}(p)} \leq 1$.

The speed-up is bounded by the number of processors and the efficiency is thus in $[0, 1]$.

Still, **Supra-linear speed-up may happen!**

Speed-up

We have $\text{SequentialTime} \leq C(p) \leq p \times \text{TotalTime}(p)$.

Hence, $s(p) = \frac{\text{SequentialTime}}{\text{TotalTime}(p)} \leq p$ and $e(p) = \frac{\text{SequentialTime}}{p \text{TotalTime}(p)} \leq 1$.

The speed-up is bounded by the number of processors and the efficiency is thus in $[0, 1]$.

Still, **Supra-linear speed-up may happen!**

We did not take I/O into account. With p processor, we have p times more available memory. Swapping sometimes kills the sequential algorithm.

Speed-up

We have $\text{SequentialTime} \leq C(p) \leq p \times \text{TotalTime}(p)$.

Hence, $s(p) = \frac{\text{SequentialTime}}{\text{TotalTime}(p)} \leq p$ and $e(p) = \frac{\text{SequentialTime}}{p \text{TotalTime}(p)} \leq 1$.

The speed-up is bounded by the number of processors and the efficiency is thus in $[0, 1]$.

Still, **Supra-linear speed-up may happen!**

We did not take I/O into account. With p processor, we have p times more available memory. Swapping sometimes kills the sequential algorithm.

Efficiency

$\text{TotalTime}(p)$ does not necessarily decrease with p due to the **parallelization overhead**.

Using more processors may hurt and may be particularly inefficient!

Parallel Matrix Algorithm

Block Version of the Outer-Product Algorithm

```
1 var A, B, C: array[0..m-1,0..m-1] of real
2 var bufferA, bufferB: array[0..m-1,0..m-1] of real
3 q ← SRQT(NUM_PROCS())
4 myrow ← MY_PROC_ROW()
5 mycol ← MY_PROC_COL()
6 for k = 0 to q - 1 do
7     for i = 0 to m - 1 do                { Broadcast A along rows }
8         BROADCASTROW(i, k, A, bufferA, m × m)
9     for j = 0 to m - 1 do                { Broadcast B along columns }
10        BROADCASTCOL(k, j, B, bufferB, m × m)
11                                           { Multiply matrix blocks }
12    if (myrow = k) And (mycol = k) then
13        MATRIXMULTIPLYADD(C, A, B, m)
14    else if (myrow = k) then MATRIXMULTIPLYADD(C, bufferA, B, m)
15    else if (mycol = k) then MATRIXMULTIPLYADD(C, A, bufferB, m)
16    else MATRIXMULTIPLYADD(C, bufferA, bufferB, m)
```

Parallel Matrix Algorithm

Block Version of the Outer-Product Algorithm

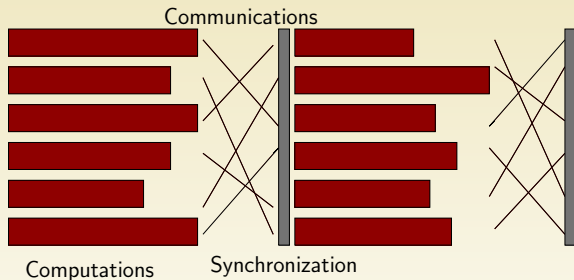
```
1 var A, B, C: array[0..m-1,0..m-1] of real
2 var bufferA, bufferB: array[0..m-1,0..m-1] of real
3 q ← SRQT(NUM_PROCS())
4 myrow ← MY_PROC_ROW()
```

Two Comments

- ▶ Many parallel programs take the number of processors as an input and adapt to it.
- ▶ Many parallel programs use collective communication operations and synchronization.

```
11 { Multiply matrix blocks }
12 if (myrow = k) And (mycol = k) then
13   | MATRIXMULTIPLYADD(C, A, B, m)
14   else if (myrow = k) then MATRIXMULTIPLYADD(C, bufferA, B, m)
15   else if (mycol = k) then MATRIXMULTIPLYADD(C, A, bufferB, m)
16   else MATRIXMULTIPLYADD(C, bufferA, bufferB, m)
```

Bulk Synchronous Parallel is a programming paradigm whose principle is a series of independent steps of computations and communication/synchronization.



The cost of a superstep is determined as the sum of three terms:

$$T = \max_i w(i) + \max h(i)g + l$$

Scheduling under BSP is about finding a tradeoff between load-balancing and number of communication/synchronizations.

- 1 **Modeling Applications, General Notions**
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 **Defining a Scheduling Problem**
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 **Batch Scheduling**
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 **Gang Scheduling as an Alternative**
 - Principles

Moldable Parallel Programs

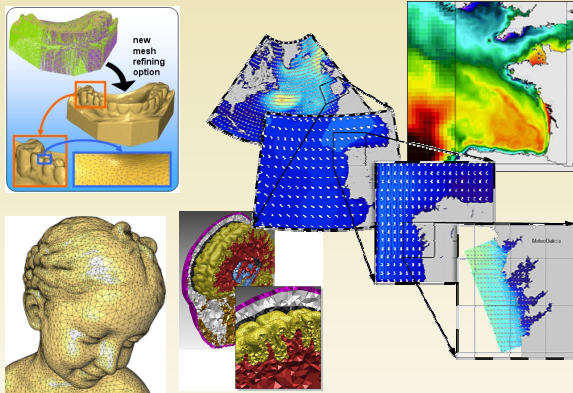
Remember the previous “Block Version of the Outer-Product Algorithm”.

```
1  $q \leftarrow \text{SRQT}(\text{NUM\_PROCS}())$   
2  $myrow \leftarrow \text{MY\_PROC\_ROW}()$   
3  $mycol \leftarrow \text{MY\_PROC\_COL}()$   
4 for  $k = 0$  to  $q - 1$  do  
5   | ...
```

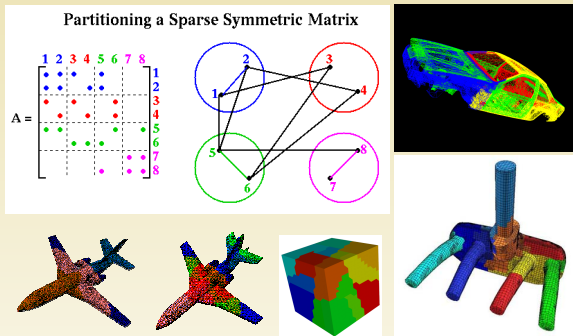
This q is not hard-coded. The algorithm adapts to the **number of available processors at the beginning** of the execution. It uses this number to distribute the data and organize the communications.

Such programs are called **moldable**.

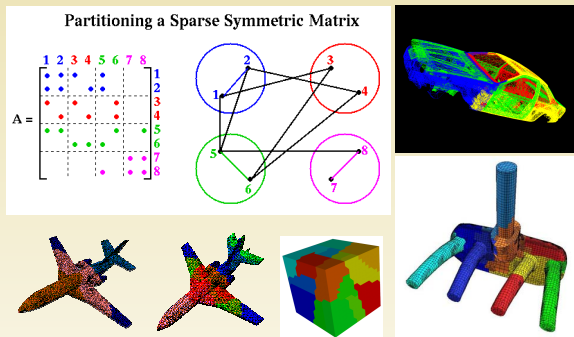
Code Coupling and Adaptive Mesh-Refining



Mesh Partitioning

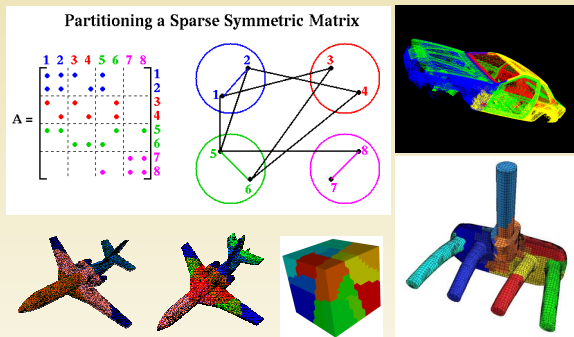


Code Coupling and Adaptive Mesh-Refining



When using adaptive mesh-refining, **load imbalance** occurs. Coupling code makes it worse.

Code Coupling and Adaptive Mesh-Refining

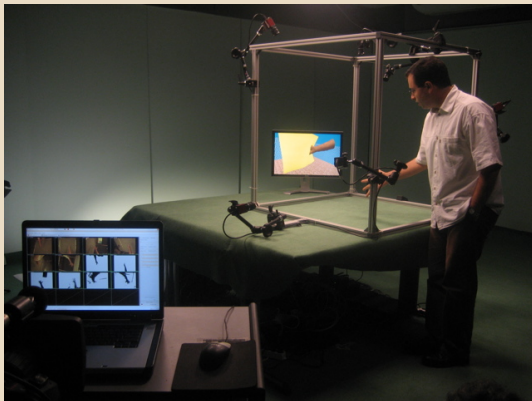


When using adaptive mesh-refining, **load imbalance** occurs. Coupling code makes it worse. Recomputing a good partition and redistributing the data is not necessarily a good option. However **adding computing resources** on the fly is often very efficient.

↪ The resource requirements vary over the time.

This kind of program is called **dynamic**.

FlowVR: Adaptive Interactive Rendering



These programs can **adapt** to the resource they are allotted over the time.

This kind of program is called **malleable**.

KAAPI: Adaptive, Asynchronous Parallel and Interactive Computing

KAAPI is based on **work-stealing** algorithms and contains **non-blocking** and **scalable** algorithms.

KAAPI/Taktuk won the 4th and 5th International Challenge GRIDS@WORK (2007, 2008).

2007 N-queens

2008 Super Quant Monte-Carlo, pricing application.

- ▶ 3609 cores used between France and Japan during one hour.
- ▶ The KAAPI/Taktuk team was able to price 988 actions on the 1000 of the challenge and was scored 8760/18000.
- ▶ The second team was able to price 177 actions using 4329 and was scored 1459/18000.

These programs can **adapt** to the resource they are allotted over the time.

This kind of program is called **malleable**.

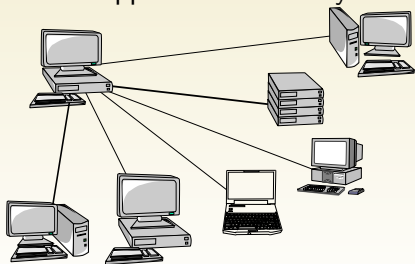
Divisible Load Scheduling

Parallelizing generally has a price. There is a **computation overhead** and a **communication/synchronization overhead**.

Some applications however have a very low computation overhead and can be very easily divided.

- ▶ Pattern Searching
- ▶ Database Computation
- ▶ Video encoding
- ▶ Image processing

Such applications are very well suited to **master-slave** computing.



Computation times and communication times are **linear** with the fraction of allotted load.

This kind of program is called **divisible**.

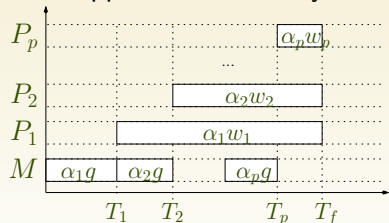
Divisible Load Scheduling

Parallelizing generally has a price. There is a **computation overhead** and a **communication/synchronization overhead**.

Some applications however have a very low computation overhead and can be very easily divided.

- ▶ Pattern Searching
- ▶ Database Computation
- ▶ Video encoding
- ▶ Image processing

Such applications are very well suited to **master-slave** computing.



Computation times and communication times are **linear** with the fraction of allotted load.

This kind of program is called **divisible**.

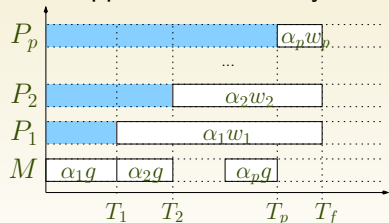
Divisible Load Scheduling

Parallelizing generally has a price. There is a **computation overhead** and a **communication/synchronization overhead**.

Some applications however have a very low computation overhead and can be very easily divided.

- ▶ Pattern Searching
- ▶ Database Computation
- ▶ Video encoding
- ▶ Image processing

Such applications are very well suited to **master-slave** computing.



Computation times and communication times are **linear** with the fraction of allotted load.

This kind of program is called **divisible**.

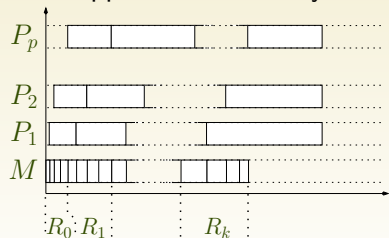
Divisible Load Scheduling

Parallelizing generally has a price. There is a **computation overhead** and a **communication/synchronization overhead**.

Some applications however have a very low computation overhead and can be very easily divided.

- ▶ Pattern Searching
- ▶ Database Computation
- ▶ Video encoding
- ▶ Image processing

Such applications are very well suited to **master-slave** computing.



Computation times and communication times are **linear** with the fraction of allotted load.

This kind of program is called **divisible**.

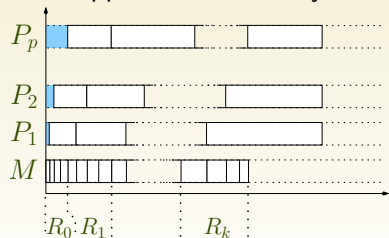
Divisible Load Scheduling

Parallelizing generally has a price. There is a **computation overhead** and a **communication/synchronization overhead**.

Some applications however have a very low computation overhead and can be very easily divided.

- ▶ Pattern Searching
- ▶ Database Computation
- ▶ Video encoding
- ▶ Image processing

Such applications are very well suited to **master-slave** computing.



Computation times and communication times are **linear** with the fraction of allotted load.

This kind of program is called **divisible**.

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

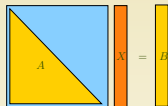
Analyzing a Simple Code

Solving $A.x = B$ where A is lower triangular matrix
for $i = 1$ **to** n **do**

Task $T_{i,i}$: $x(i) \leftarrow b(i)/a(i,i)$

for $j = i + 1$ **to** n **do**

Task $T_{i,j}$: $b(j) \leftarrow b(j) - a(j,i) \times x(i)$



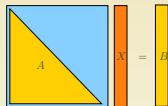
Analyzing a Simple Code

Solving $A.x = B$ where A is lower triangular matrix
for $i = 1$ **to** n **do**

Task $T_{i,i}$: $x(i) \leftarrow b(i)/a(i,i)$

for $j = i + 1$ **to** n **do**

Task $T_{i,j}$: $b(j) \leftarrow b(j) - a(j,i) \times x(i)$



For a given value $1 \leq i \leq n$, all tasks $T_{i,*}$ are computations done during the i^{th} iteration of the outer loop.

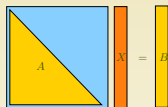
Analyzing a Simple Code

Solving $A.x = B$ where A is lower triangular matrix
for $i = 1$ **to** n **do**

Task $T_{i,i}$: $x(i) \leftarrow b(i)/a(i,i)$

for $j = i + 1$ **to** n **do**

Task $T_{i,j}$: $b(j) \leftarrow b(j) - a(j,i) \times x(i)$



For a given value $1 \leq i \leq n$, all tasks $T_{i,*}$ are computations done during the i^{th} iteration of the outer loop.

\prec_{seq} is the **sequential order** :

$$T_{1,1} \prec_{seq} T_{1,2} \prec_{seq} T_{1,3} \prec_{seq} \dots \prec_{seq} T_{1,n} \prec_{seq} T_{2,2} \prec_{seq} T_{2,3} \prec_{seq} \dots \prec_{seq} T_{n,n} .$$

However, some **independent** tasks could be executed in parallel. Independent tasks are the ones whose execution order can be changed without modifying the result of the program. Two independent tasks may read the value but never write to the same memory location.

However, some **independent** tasks could be executed in parallel. Independent tasks are the ones whose execution order can be changed without modifying the result of the program. Two independent tasks may read the value but never write to the same memory location.

For a given task T , $In(T)$ denotes the set of input variables and $Out(T)$ the set of output variables.

However, some **independent** tasks could be executed in parallel. Independent tasks are the ones whose execution order can be changed without modifying the result of the program. Two independent tasks may read the value but never write to the same memory location.

For a given task T , $In(T)$ denotes the set of input variables and $Out(T)$ the set of output variables.

In the previous example, we have :

$\left\{ \begin{array}{l} In(T_{i,i}) = \{b(i), a(i, i)\} \\ Out(T_{i,i}) = \{x(i)\} \text{ and} \\ In(T_{i,j}) = \{b(j), a(j, i), x(i)\} \\ Out(T_{i,j}) = \{b(j)\} \text{ for } j > i. \end{array} \right.$	for $i = 1$ to n do									
	<table border="0"><tr><td>Task $T_{i,i}$:</td><td>$x(i) \leftarrow b(i)/a(i, i)$</td></tr><tr><td>for $j = i + 1$ to n do</td><td></td></tr><tr><td><table border="0"><tr><td>Task $T_{i,j}$:</td><td>$b(j) \leftarrow b(j) - a(j, i) \times$</td></tr><tr><td></td><td>$x(i)$</td></tr></table></td><td></td></tr></table>	Task $T_{i,i}$:	$x(i) \leftarrow b(i)/a(i, i)$	for $j = i + 1$ to n do		<table border="0"><tr><td>Task $T_{i,j}$:</td><td>$b(j) \leftarrow b(j) - a(j, i) \times$</td></tr><tr><td></td><td>$x(i)$</td></tr></table>	Task $T_{i,j}$:	$b(j) \leftarrow b(j) - a(j, i) \times$		$x(i)$
Task $T_{i,i}$:	$x(i) \leftarrow b(i)/a(i, i)$									
for $j = i + 1$ to n do										
<table border="0"><tr><td>Task $T_{i,j}$:</td><td>$b(j) \leftarrow b(j) - a(j, i) \times$</td></tr><tr><td></td><td>$x(i)$</td></tr></table>	Task $T_{i,j}$:	$b(j) \leftarrow b(j) - a(j, i) \times$		$x(i)$						
Task $T_{i,j}$:	$b(j) \leftarrow b(j) - a(j, i) \times$									
	$x(i)$									

Definition.

Two tasks T and T' are not independent ($T \perp T'$) whenever they share a written variable:

$$T \perp T' \Leftrightarrow \left\{ \begin{array}{l} In(T) \cap Out(T') \neq \emptyset \\ \text{or } Out(T) \cap In(T') \neq \emptyset \\ \text{or } Out(T) \cap Out(T') \neq \emptyset \end{array} \right. .$$

Those conditions are known as Bernstein's conditions [Bernstein66].

Definition.

Two tasks T and T' are not independent ($T \perp T'$) whenever they share a written variable:

$$T \perp T' \Leftrightarrow \begin{cases} In(T) \cap Out(T') \neq \emptyset \\ \text{or } Out(T) \cap In(T') \neq \emptyset \\ \text{or } Out(T) \cap Out(T') \neq \emptyset \end{cases} .$$

Those conditions are known as Bernstein's conditions [Bernstein66].

We can check that:

```
for  $i = 1$  to  $n$  do
  Task  $T_{i,i}$ :  $x(i) \leftarrow b(i)/a(i, i)$ 
  for  $j = i + 1$  to  $n$  do
    Task  $T_{i,j}$ :  $b(j) \leftarrow b(j) - a(j, i) \times x(i)$ 
```

Definition.

Two tasks T and T' are not independent ($T \perp T'$) whenever they share a written variable:

$$T \perp T' \Leftrightarrow \begin{cases} In(T) \cap Out(T') \neq \emptyset \\ \text{or } Out(T) \cap In(T') \neq \emptyset \\ \text{or } Out(T) \cap Out(T') \neq \emptyset \end{cases} .$$

Those conditions are known as Bernstein's conditions [Bernstein66].

We can check that:

$$\begin{aligned} \blacktriangleright Out(T_{1,1}) \cap In(T_{1,2}) &= \{x(1)\} \text{ for } i = 1 \text{ to } n \text{ do} \\ &\quad \left[\begin{array}{l} \text{Task } T_{i,i}: x(i) \leftarrow b(i)/a(i, i) \\ \text{for } j = i + 1 \text{ to } n \text{ do} \\ \quad \left[\text{Task } T_{i,j}: b(j) \leftarrow b(j) - a(j, i) \times x(i) \right] \end{array} \right. \\ &\rightsquigarrow T_{1,1} \perp T_{1,2}. \end{aligned}$$

Definition.

Two tasks T and T' are not independent ($T \perp T'$) whenever they share a written variable:

$$T \perp T' \Leftrightarrow \begin{cases} In(T) \cap Out(T') \neq \emptyset \\ \text{or } Out(T) \cap In(T') \neq \emptyset \\ \text{or } Out(T) \cap Out(T') \neq \emptyset \end{cases} .$$

Those conditions are known as Bernstein's conditions [Bernstein66].

We can check that:

- ▶ $Out(T_{1,1}) \cap In(T_{1,2}) = \{x(1)\}$ for $i = 1$ to n do
 - Task $T_{i,i}$: $x(i) \leftarrow b(i)/a(i, i)$ $\leadsto T_{1,1} \perp T_{1,2}$.
- ▶ $Out(T_{1,3}) \cap Out(T_{2,3}) = \{b(3)\}$ for $j = i + 1$ to n do
 - Task $T_{i,j}$: $b(j) \leftarrow b(j) - a(j, i) \times x(i)$ $\leadsto T_{1,3} \perp T_{2,3}$.

Precedence

If $T \perp T'$, then they should be ordered with the sequential execution order. $T \prec T'$ if $T \perp T'$ and $T <_{seq} T'$.

More precisely \prec is defined as the **transitive closure** of $(<_{seq} \cap \perp)$.

Precedence

If $T \perp T'$, then they should be ordered with the sequential execution order. $T \prec T'$ if $T \perp T'$ and $T <_{seq} T'$.

More precisely \prec is defined as the **transitive closure** of $(<_{seq} \cap \perp)$.

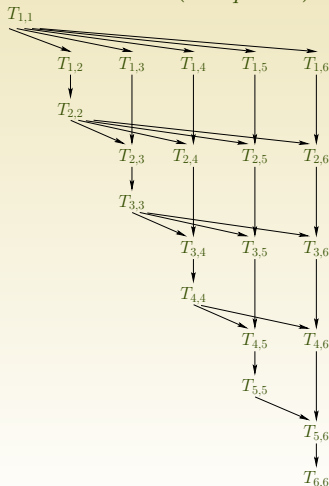
for $i = 1$ to n do

Task $T_{i,i}$: $x(i) \leftarrow b(i)/a(i, i)$

for $j = i + 1$ to n do

Task $T_{i,j}$: $b(j) \leftarrow b(j) - a(j, i) \times x(i)$

A **dependence graph** G is used.



Precedence

If $T \perp T'$, then they should be ordered with the sequential execution order. $T \prec T'$ if $T \perp T'$ and $T <_{seq} T'$.

More precisely \prec is defined as the **transitive closure** of $(<_{seq} \cap \perp)$.

for $i = 1$ to n do

Task $T_{i,i}$: $x(i) \leftarrow b(i)/a(i, i)$

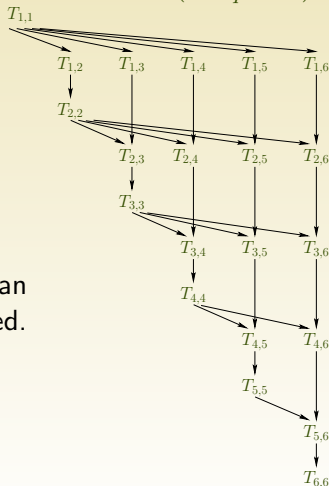
for $j = i + 1$ to n do

Task $T_{i,j}$: $b(j) \leftarrow b(j) - a(j, i) \times x(i)$

A **dependence graph** G is used.

$(e : T \rightarrow T') \in G$ means that T' can start only if T has already been finished.

T is a **predecessor** of T' .



Precedence

If $T \perp T'$, then they should be ordered with the sequential execution order. $T \prec T'$ if $T \perp T'$ and $T <_{seq} T'$.

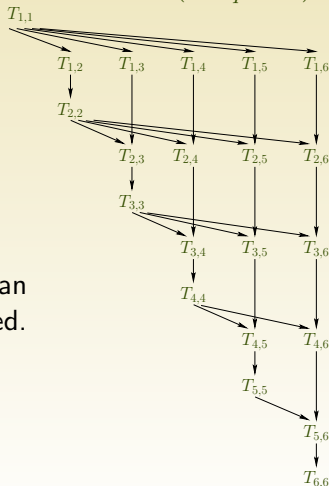
More precisely \prec is defined as the **transitive closure** of $(<_{seq} \cap \perp)$.

for $i = 1$ to n do

Task $T_{i,i}$: $x(i) \leftarrow b(i)/a(i, i)$

for $j = i + 1$ to n do

Task $T_{i,j}$: $b(j) \leftarrow b(j) - a(j, i) \times x(i)$



A **dependence graph** G is used.

$(e : T \rightarrow T') \in G$ means that T' can start only if T has already been finished.

T is a **predecessor** of T' .

Transitivity arcs are generally omitted.

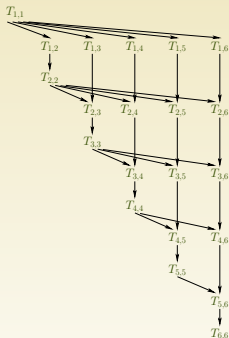
The previous task graph comes from a **low-level** analysis of the code.

It probably makes little sense to do a parallel implementation with **MPI** with such a low task granularity.

Can totally make sense with **OpenMP**.

Such task graphs can also be used by compilers to do code optimization by exploiting multiple functional units, pipelines functional units, etc.

With **blocking** these tasks could become MPI (parallel) tasks.



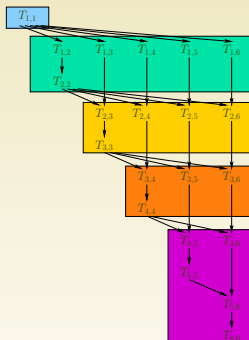
The previous task graph comes from a **low-level** analysis of the code.

It probably makes little sense to do a parallel implementation with **MPI** with such a low task granularity.

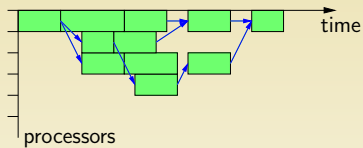
Can totally make sense with **OpenMP**.

Such task graphs can also be used by compilers to do code optimization by exploiting multiple functional units, pipelines functional units, etc.

With **blocking** these tasks could become MPI (parallel) tasks.

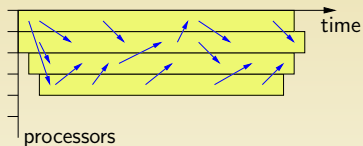


... to Parallel Tasks

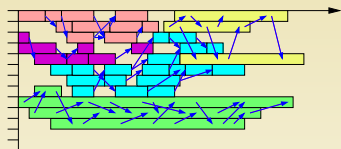


Hide applications' complexity

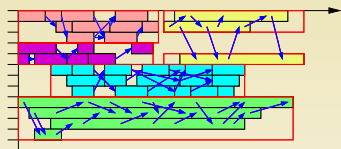
... to Parallel Tasks



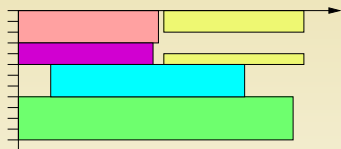
Hide applications' complexity



Hide applications' complexity



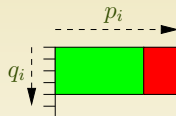
Hide applications' complexity



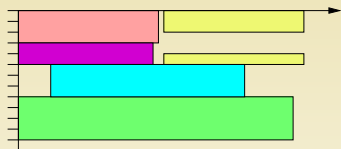
Hide applications' complexity

3 versions:

- ▶ Rigid Tasks



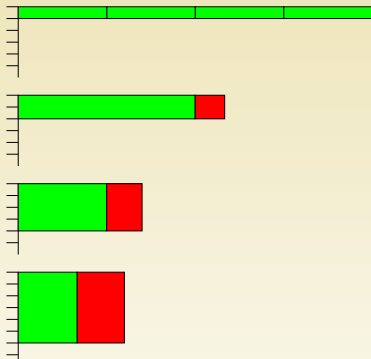
The execution time *generally* decreases with the number of processors but the penalty incurred by communications and synchronizations increases.



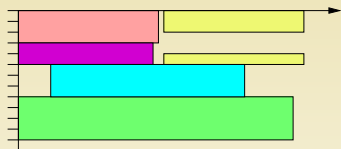
Hide applications' complexity

3 versions:

- ▶ Rigid Tasks
- ▶ Moldable Tasks



The execution time *generally* decreases with the number of processors but the penalty incurred by communications and synchronizations increases.



Hide applications' complexity

3 versions:

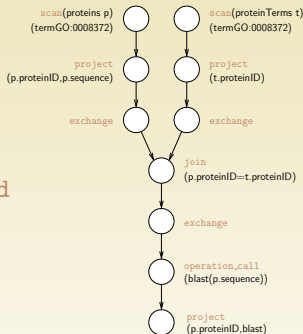
- ▶ Rigid Tasks
- ▶ Moldable Tasks
- ▶ Malleable Tasks



The execution time *generally* decreases with the number of processors but the penalty incurred by communications and synchronizations increases.

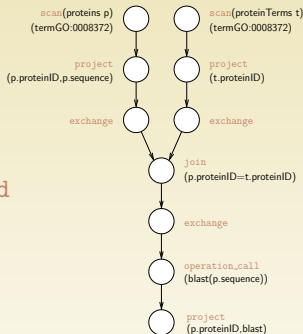
Task-graph do not necessarily come from instruction-level analysis.

```
select p.proteinID,  
       blast(p.sequence)  
from proteins p, proteinTerms t  
where p.proteinID = t.proteinID and  
t.term = GO:0008372
```



Task-graph do not necessarily come from instruction-level analysis.

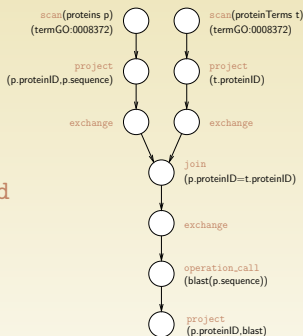
```
select p.proteinID,  
       blast(p.sequence)  
from proteins p, proteinTerms t  
where p.proteinID = t.proteinID and  
t.term = GO:0008372
```



- ▶ Each task may be a parallel job. . .

Task-graph do not necessarily come from instruction-level analysis.

```
select p.proteinID,  
       blast(p.sequence)  
from proteins p, proteinTerms t  
where p.proteinID = t.proteinID and  
t.term = GO:0008372
```



- ▶ Each task may be a parallel job. . .
- ▶ Each edge depicts a dependency i.e. most of the times some data to transfer.

I have presented you a few different parallel program models:

- ▶ rigid jobs
- ▶ moldable jobs
- ▶ dynamic jobs
- ▶ malleable jobs
- ▶ divisible jobs
- ▶ BSP jobs
- ▶ DAGs of the previous jobs

The rationale behind all these models is:

- ▶ the diversity and the complexity of parallel programs;
- ▶ the level of details we need/wish to expose to the one in charge of the execution.

Modeling is an art.

You have to know your application to know what is negligible and what is important. Even if your model is imperfect, you may still derive interesting results.

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Some questions you need to answer

Preemption Are we allowed to suspend a program and resume it later ?

- ▶ Resumed from the beginning or from where it was stopped?
- ▶ May be resumed on another machine or not (migration)?
- ▶ Does preemption/migration has a cost or not?

Release dates Are all tasks available at the very beginning or not?

Deadlines Are the tasks associated to a deadline before which they should complete? What happens when the deadline is missed?

Dependencies Are there dependencies between tasks (DAGs)?

Users Are there many users and should this be taken into account?

Long-term vs. short-term What kind of constraints do you have on the time needed to take your scheduling decisions?

Online vs. Off-line, clairvoyance

What kind of information do you have to make your scheduling choices?

Off-line You know everything (release dates and processing time of each task) at the very beginning.

It is the “simplest” setting and will give you insights on your scheduling problem even though these hypothesis do not really hold in practice.

This kind of problem should thus be studied *before everything else*.

On-line/clairvoyant You do not know in advance when tasks arrive. However, once a new task are available, you know its computation time.

On-line/non clairvoyant You know nothing!

Sometimes (often?), reality is in between:

- ▶ We could have “informations” about the task arrival (e.g., periodic creation, random process, use the past to predict the future).
- ▶ We could have “informations” about the task computation requirement (e.g., mix of short tasks and long tasks).

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - **Criteria: How Do You Win the Game?**
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Criteria: Intuitive Notion

CPU utilization (max) percent usage of CPU. Only useful computations (mix CPU, I/O; preemption overhead).

Throughput (max) *average* number of tasks that complete their execution per time-unit.

Makespan (min) Completion time of the last finishing task.

Load (min) Completion time of the last finishing task for a given processor.

Turnaround Time/Response Time/Flow (min) amount of time it takes between the task arrival and its completion.

Waiting Time (min) amount of time spent waiting for being executed.

Slowdown/Stretch (min) slowdown factor encountered by a task relative to the time it would take on an unloaded system.

The previous quantities are task- or CPU-centric and need to be aggregated into a single objective function.

- ▶ max (the worst case)
- ▶ average: arithmetic (i.e. sum) or something else...
- ▶ variance (to be “fair” between the tasks).

A given task T_i is defined by:

- ▶ processing time p_i
- ▶ release date r_i
- ▶ completion time C_i
- ▶ (number of required processors q_i)
- ▶ (deadline d_i)

Completion Time

- ▶ Makespan: $C_{\max} = \max_i C_i$
This metric is the most classical and is relevant when scheduling a *single* application.
- ▶ Total Completion Time: $SC = \sum_i C_i$

A given task T_i is defined by:

- ▶ processing time p_i
- ▶ release date r_i
- ▶ completion time C_i
- ▶ (number of required processors q_i)
- ▶ (deadline d_i)

Response Time

$$F_i = C_i - r_i$$

- ▶ Maximum Flow Time: $F_{\max} = \max_i F_i$
- ▶ Total Completion Time: $SF = \sum_i F_i = SC - \sum_i r_i$

A given task T_i is defined by:

- ▶ processing time p_i
- ▶ release date r_i
- ▶ completion time C_i
- ▶ (number of required processors q_i)
- ▶ (deadline d_i)

Waiting Time

$$W_i = C_i - r_i - p_i$$

- ▶ Maximum Waiting time: $W_{\max} = \max_i W_i$
- ▶ Total Waiting Time: $SW = \sum_i W_i = SF - \sum_i p_i$

A given task T_i is defined by:

- ▶ processing time p_i
- ▶ release date r_i
- ▶ completion time C_i
- ▶ (number of required processors q_i)
- ▶ (deadline d_i)

Slowdown

$$S_i = \frac{C_i - r_i}{p_i}$$

- ▶ Maximum Stretch: $S_{\max} = \max_i S_i$
- ▶ Total Stretch: $SS = \sum_i S_i$

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - **Analysis Method**
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Most scheduling problem are NP-complete but you may be lucky. . . So the first question to answer is: P or NP-hard ?

For a given objective function Obj :

Definition: Decision and Optimization.

$Dec(M)$: Is there a schedule σ such that $Obj(\sigma) \leq M$?

Opt : Find M^* such that $M^* = \min_{\sigma} Obj(\sigma)$.

If Dec can be solved in polynomial time, then so can Opt (using a dichotomy). And conversely. . .

Note that since $SW(\sigma) = SF(\sigma) - \sum_i p_i = SC(\sigma) - \sum_i r_i - \sum_i p_i$, all these problem are equivalent on a complexity point of view.

Your scheduling problem is NP-hard so you need to propose a heuristic and compare it to the best possible solution.
Consider a given objective function Obj .

Definition: ϱ -approximation.

An algorithm \mathcal{A} is a ϱ -approximation
iff
for any instance I , $Obj(\mathcal{A}(I)) \leq \varrho \cdot Obj^*(I)$.

The approximation ratio of \mathcal{A} is:

$$\varrho(\mathcal{A}) = \max_I \frac{Obj(\mathcal{A}(I))}{Obj^*(I)}$$

Note that even though $SW(\sigma) = SF(\sigma) - \sum_i p_i = SC(\sigma) - \sum_i r_i - \sum_i p_i$, these problem are **not** equivalent on an approximation point of view.

What is the best solution to an online problem (where the heuristic doesn't know in advance the jobs arrival) ?

We keep comparing to the best possible solution, i.e. the one that knows everything.

Definition: ρ -competitive.

An algorithm \mathcal{A} is a ρ -approximation
iff
for any instance I , $Obj(\mathcal{A}(I)) \leq \rho \cdot Obj^*(I)$.

The approximation ratio of \mathcal{A} is:

$$\rho(\mathcal{A}) = \max_I \frac{Obj(\mathcal{A}(I))}{Obj^*(I)}$$

It is the same definition except that it applies to online algorithms.
For such a pessimistic evaluation, one commonly uses an **adversary**.

Average-Case Analysis

If we have a probability distribution over the set of instances, Obj can thus be seen as a random variable.

We can define the expectation of Obj .

$$\mathbb{E}[Obj(\mathcal{A})] = \int_I Obj(\mathcal{A}(I))p(I).dI = \sum_I Obj(\mathcal{A}(I))p(I)$$

People often try to evaluate at (at least through experiments)

$$\varrho(\mathcal{A}) = \int_I \frac{Obj(\mathcal{A}(I))}{Obj^*(I)}p(I).dI \text{ or}$$
$$\varrho(\mathcal{A}) = \frac{\mathbb{E}[Obj(\mathcal{A})]}{\mathbb{E}[Obj^*]} = \frac{\int_I Obj(\mathcal{A}(I))p(I).dI}{\int_I Obj^*(I)p(I).dI}$$

However, in the literature, there are many different ways of comparing random variables (and thus to compare and evaluate algorithms). These techniques will be presented in much more details later.

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - **Graham Notation**
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Graham Notation

Many parameter can change in a scheduling problem. Graham has then proposed the following classification : $\langle \alpha | \beta | \gamma \rangle$ [Brucker-Book]

Graham Notation

Many parameter can change in a scheduling problem. Graham has then proposed the following classification : $\langle \alpha | \beta | \gamma \rangle$ [Brucker-Book]

- ▶ α is the processor environment (a few examples):
 - ▶ \emptyset : single processor;
 - ▶ P : identical processors;
 - ▶ Q : uniform processors;
 - ▶ R : unrelated processors;

Graham Notation

Many parameter can change in a scheduling problem. Graham has then proposed the following classification : $\langle \alpha | \beta | \gamma \rangle$ [Brucker-Book]

- ▶ α is the processor environment (a few examples):
 - ▶ \emptyset : single processor;
 - ▶ P : identical processors;
 - ▶ Q : uniform processors;
 - ▶ R : unrelated processors;
- ▶ β describe task and resource characteristics (a few examples):
 - ▶ *pmtn*: preemption;
 - ▶ *prec*, *tree* or *chains*: general precedence constraints, tree constraints, set of chain constraints and independent tasks otherwise;
 - ▶ r_j : tasks have release dates;
 - ▶ $p_j = p$ or $\underline{p} \leq p_j \leq \bar{p}$: all task have processing time equal to p , or comprised between \underline{p} and \bar{p} , or have arbitrary processing times otherwise;
 - ▶ \tilde{d} : deadlines;

Graham Notation

Many parameter can change in a scheduling problem. Graham has then proposed the following classification : $\langle \alpha | \beta | \gamma \rangle$ [Brucker-Book]

- ▶ α is the processor environment (a few examples):
 - ▶ \emptyset : single processor;
 - ▶ P : identical processors;
 - ▶ Q : uniform processors;
 - ▶ R : unrelated processors;
- ▶ β describe task and resource characteristics (a few examples):
 - ▶ *pmtn*: preemption;
 - ▶ *prec*, *tree* or *chains*: general precedence constraints, tree constraints, set of chain constraints and independent tasks otherwise;
 - ▶ r_j : tasks have release dates;
 - ▶ $p_j = p$ or $\underline{p} \leq p_j \leq \bar{p}$: all task have processing time equal to p , or comprised between \underline{p} and \bar{p} , or have arbitrary processing times otherwise;
 - ▶ \tilde{d} : deadlines;
- ▶ γ denotes the optimization criterion (a few examples):
 - ▶ C_{\max} : makespan;
 - ▶ L_{\max} : maximum lateness ($\max C_i - d_i$);
 - ▶ $\sum C_i$: average completion time;
 - ▶ \dots
 - ▶ $\sum w_i C_i$: weighted A.C.T;

Understand the following problems and propose a practical situation to illustrate them:

- ▶ $\langle P|prec|C_{\max}\rangle$
- ▶ $\langle P|q_j, prec|C_{\max}\rangle$
- ▶ $\langle P|q_j|F_{\max}\rangle$
- ▶ $\langle 1|r_j; pmtn|S_{\max}\rangle$
- ▶ $\langle 1|r_j; pmtn, d_i|L_{\max}\rangle$

Scheduling is a very generic word that encompasses a very wide range of situations, problems and analysis techniques.

Scheduling is generally about deciding *who*, *where* and *when*.

It is thus almost everywhere and when you start looking at a given scheduling problem, with very high probability, many people already worked on it.

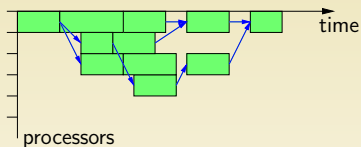
Doing a serious and **thorough bibliographical study** is thus of uttermost importance!

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Need for Batch Scheduling

- ▶ Parallel Tasks from Scientific Computations (simulation, medical)



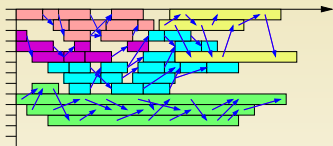
Need for Batch Scheduling

- ▶ Parallel Tasks from Scientific Computations (simulation, medical)



Need for Batch Scheduling

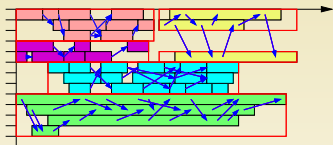
- ▶ Parallel Tasks from Scientific Computations (simulation, medical)



- ▶ When one purchases a cluster, typically **many users** want to use it.
 - ▶ One cannot let them step on each other's toes
 - ▶ Every user wants to be on a **dedicated** machine
 - ▶ Applications are written assuming some amount of RAM, some notion that all processors go at the same speed, etc.

Need for Batch Scheduling

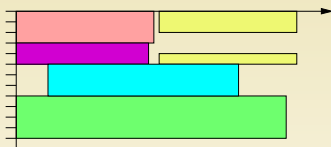
- ▶ Parallel Tasks from Scientific Computations (simulation, medical)



- ▶ When one purchases a cluster, typically **many users** want to use it.
 - ▶ One cannot let them step on each other's toes
 - ▶ Every user wants to be on a **dedicated** machine
 - ▶ Applications are written assuming some amount of RAM, some notion that all processors go at the same speed, etc.

Need for Batch Scheduling

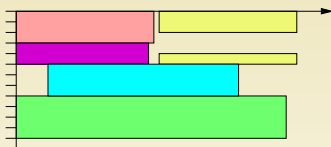
- ▶ Parallel Tasks from Scientific Computations (simulation, medical)



- ▶ When one purchases a cluster, typically **many users** want to use it.
 - ▶ One cannot let them step on each other's toes
 - ▶ Every user wants to be on a **dedicated** machine
 - ▶ Applications are written assuming some amount of RAM, some notion that all processors go at the same speed, etc.

Need for Batch Scheduling

- ▶ Parallel Tasks from Scientific Computations (simulation, medical)



- ▶ When one purchases a cluster, typically **many users** want to use it.
 - ▶ One cannot let them step on each other's toes
 - ▶ Every user wants to be on a **dedicated** machine
 - ▶ Applications are written assuming some amount of RAM, some notion that all processors go at the same speed, etc.

The **Job Scheduler** is the entity that prevents them from stepping on each other's toes

The Job Scheduler gives out nodes to applications

Batch Scheduling

Each job is defined as a **Number of nodes** (q_i) and a **Time** (p_i):

I want 6 nodes for 1h

Typically users are “charged” against an “allocation”: e.g. “*You only get 100 CPU hours per week*”.

A batch scheduler is a central middleware to manage resources (e.g. processors) of parallel machines:

- ▶ accept jobs (computing tasks) submitted by users
- ▶ decide **when** and **where** jobs are executed
- ▶ start jobs execution

They take into account:

- ▶ **unavailability** of some nodes
- ▶ users jobs **mutual exclusion**
- ▶ **specific needs** for jobs (memory, network, ...)

While trying to :

- ▶ **maximize resources usage**
- ▶ be **fair** among users

Typical wanted features:

- ▶ Interactive mode
- ▶ Batch mode
- ▶ Parallel jobs support
- ▶ Multi-queues with priorities
- ▶ Admission policies (limit on usage, notions of user groups, power users)
- ▶ Resources matching
- ▶ File staging
- ▶ Jobs dependences
- ▶ Backfilling
- ▶ Reservations
- ▶ Best effort jobs
- ▶ Environment reconfiguration

There are many existing batch schedulers : LSF, PBS/Torque, Maui scheduler, Sun Grid Engine, EASY, OAR, ...

These are **complex systems** with many config options !

Main Batch Schedulers Features

	OpenPBS	SGE	Maui Scheduler (+ OpenPBS)	OAR
Interactive mode	×	×	×	×
Batch mode	×	×	×	×
Parallel jobs support	×	×	×	×
Multi-queues with priorities	×	×	×	×
Resources matching	×	×	×	×
Admission policies	×	×	×	×
File staging	×	×	×	
Jobs dependences	×	×	×	
Backfilling			×	×
Reservations			×	×
Best effort jobs				×
Environment reconfiguration				×
Fair sharing			×	×

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 **Batch Scheduling**
 - Principles
 - **Theoretical results**
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

List Scheduling

When simple problems are hard, we should try to find good **approximation** heuristics. A ρ -approximation is an algorithm whose output is never more than a factor ρ times the optimum solution.

Natural idea: using **greedy** strategy like trying to allocate the most possible task at a given time-step. However at some point we may face a choice (when there is more ready tasks than available processors).

When simple problems are hard, we should try to find good **approximation** heuristics. A ρ -approximation is an algorithm whose output is never more than a factor ρ times the optimum solution.

Natural idea: using **greedy** strategy like trying to allocate the most possible task at a given time-step. However at some point we may face a choice (when there is more ready tasks than available processors).

Any strategy that does not let on purpose a processor idle is efficient [Coffman76]. Such a schedule is called **list-schedule**.

Theorem 1: Coffman.

Let $G = (V, E, w)$ be a DAG of sequential tasks, p the number of processors, and σ_p a list-schedule of G on p processors.

$$C_{\max}(\sigma_p) \leq \left(2 - \frac{1}{p}\right) C_{\max}^*(p).$$

When simple problems are hard, we should try to find good **approximation** heuristics. A ρ -approximation is an algorithm whose output is never more than a factor ρ times the optimum solution.

Natural idea: using **greedy** strategy like trying to allocate the most possible task at a given time-step. However at some point we may face a choice (when there is more ready tasks than available processors).

Any strategy that does not let on purpose a processor idle is efficient [Coffman76]. Such a schedule is called **list-schedule**.

Theorem 1: Coffman.

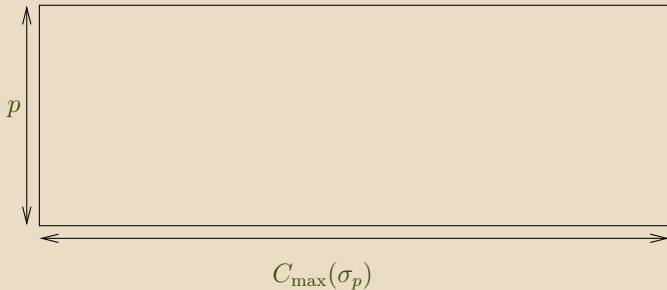
Let $G = (V, E, w)$ be a DAG of sequential tasks, p the number of processors, and σ_p a list-schedule of G on p processors.

$$C_{\max}(\sigma_p) \leq \left(2 - \frac{1}{p}\right) C_{\max}^*(p).$$

Most of the time, list-heuristics are based on the **critical path**.

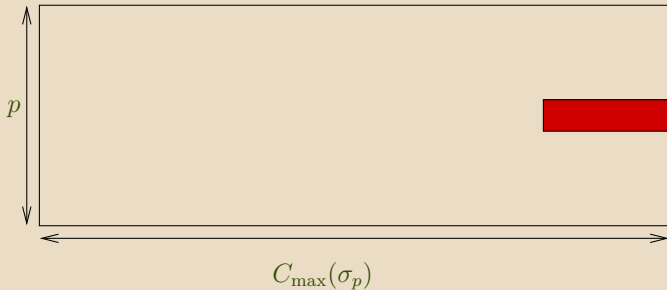
List Scheduling: proving the Coffman result

Proof.



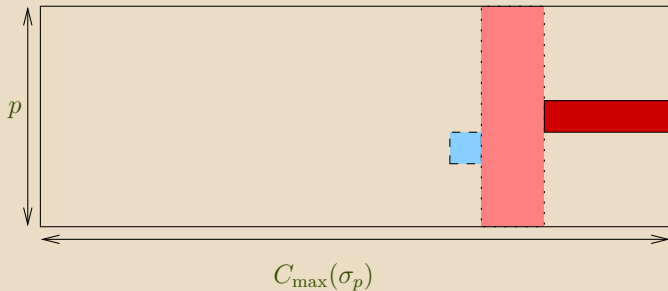
List Scheduling: proving the Coffman result

Proof.



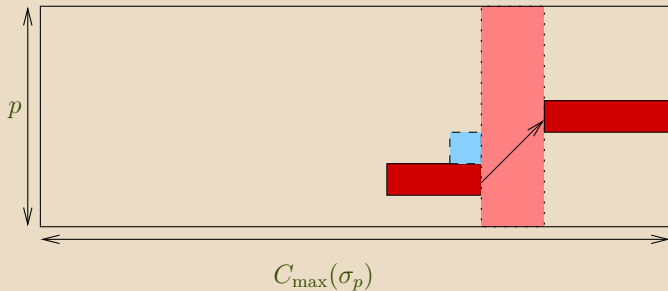
List Scheduling: proving the Coffman result

Proof.



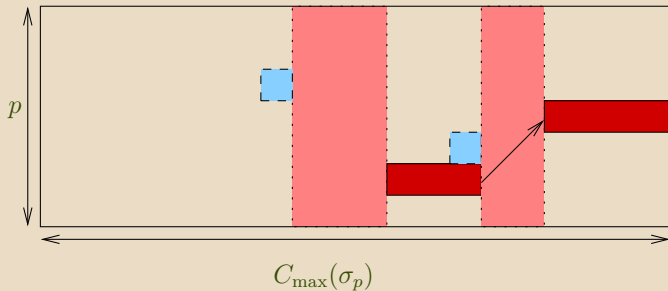
List Scheduling: proving the Coffman result

Proof.



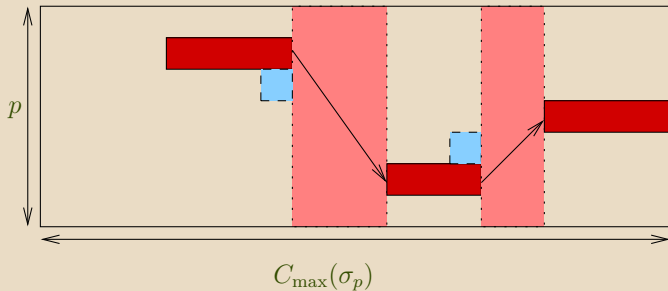
List Scheduling: proving the Coffman result

Proof.



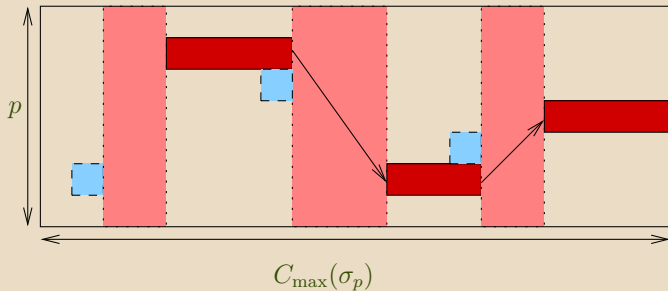
List Scheduling: proving the Coffman result

Proof.



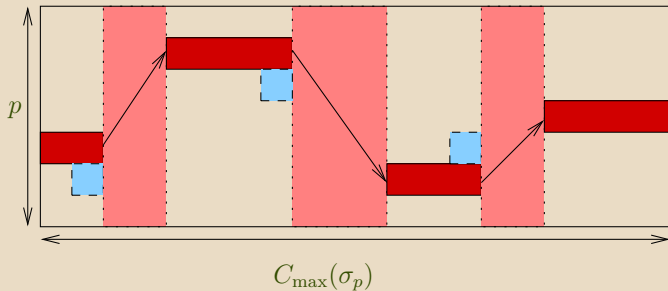
List Scheduling: proving the Coffman result

Proof.



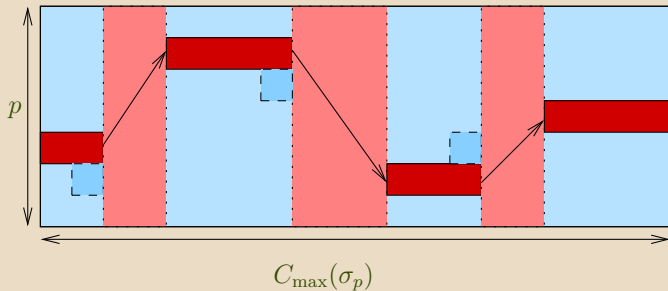
List Scheduling: proving the Coffman result

Proof.



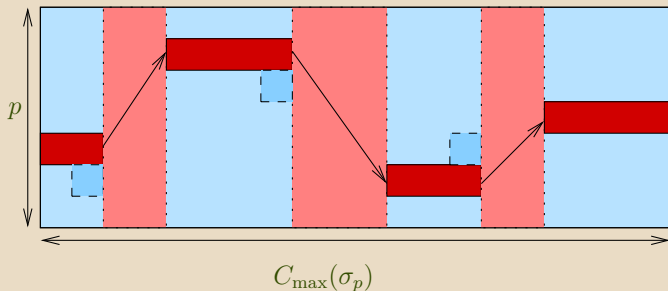
List Scheduling: proving the Coffman result

Proof.



List Scheduling: proving the Coffman result

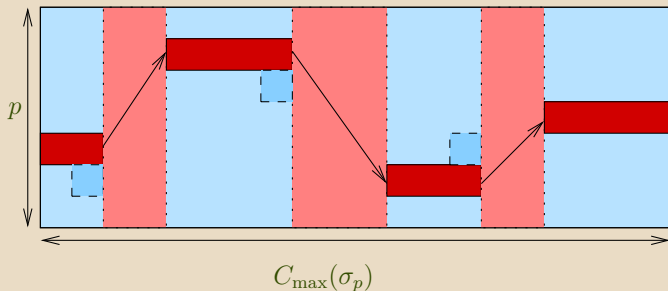
Proof.



Therefore, $Idle \leq (p - 1) \cdot w(\Phi)$ for some Φ



Proof.



Therefore, $Idle \leq (p - 1) \cdot w(\Phi)$ for some Φ

Hence,

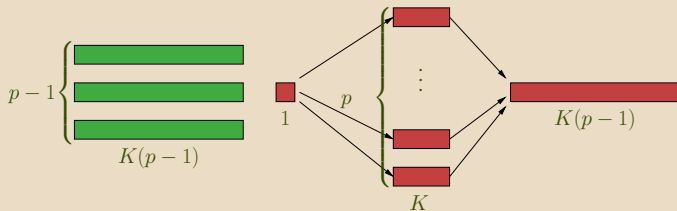
$$\begin{aligned} p \cdot C_{\max}(\sigma_p) &= Idle + Seq \leq (p - 1)w(\Phi) + Seq \\ &\leq (p - 1)C_{\max}^*(p) + p \cdot C_{\max}^*(p) = (2p - 1)C_{\max}^*(p) \end{aligned}$$

□

List Scheduling: proving the Coffman result

One can actually prove that this bound cannot be improved.

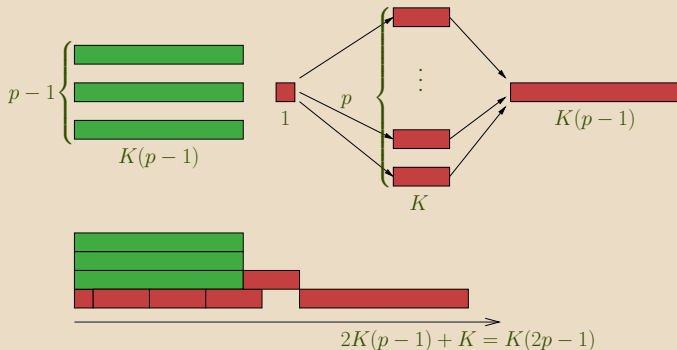
Proof.



List Scheduling: proving the Coffman result

One can actually prove that this bound cannot be improved.

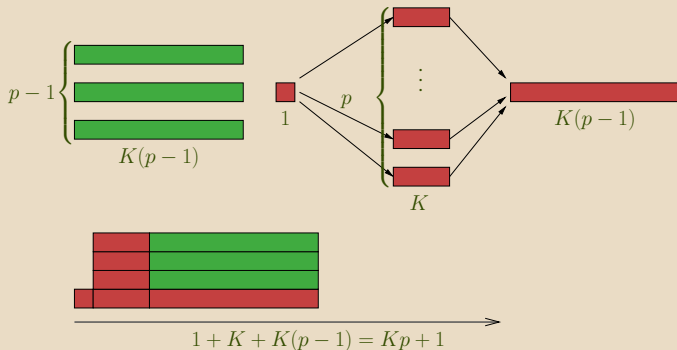
Proof.



List Scheduling: proving the Coffman result

One can actually prove that this bound cannot be improved.

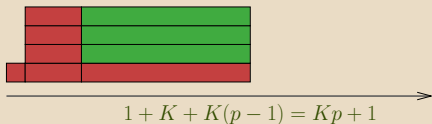
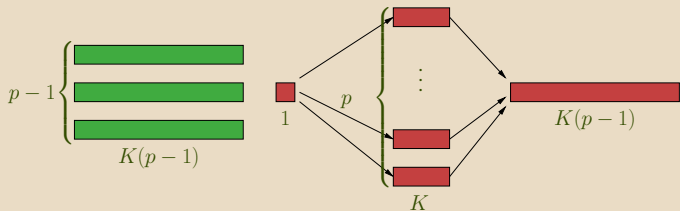
Proof.



List Scheduling: proving the Coffman result

One can actually prove that this bound cannot be improved.

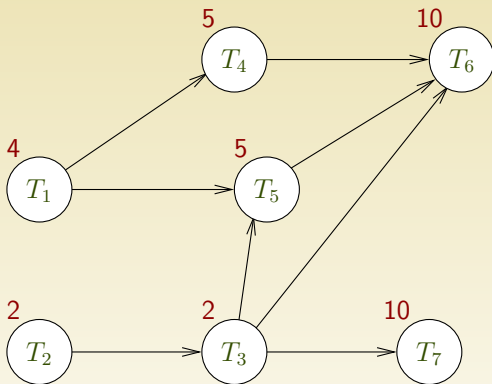
Proof.



$$\rho \geq \frac{K(2p-1)}{Kp+1} \xrightarrow{K \rightarrow \infty} \frac{2p-1}{p}$$



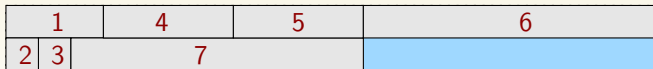
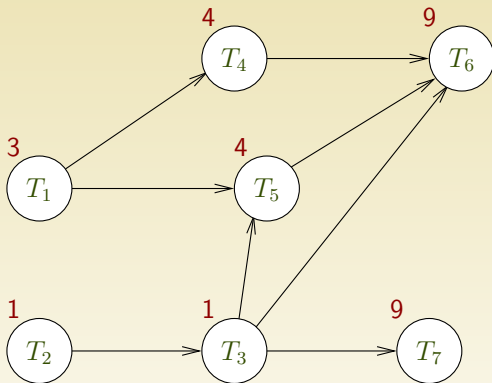
List scheduling Anomalies



1		4		6	
2	3	5		7	

$$MS = 19$$

List scheduling Anomalies



$$MS = 20$$

List Scheduling for Parallel Rigid Tasks

Let us assume we have n independent rigid jobs $J_1 = (p_1, q_1), \dots, J_n = (p_n, q_n)$ and m machines.

Let us denote by T^* the optimal makespan for this instance.

List Scheduling for Parallel Rigid Tasks

Let us assume we have n independent rigid jobs $J_1 = (p_1, q_1), \dots, J_n = (p_n, q_n)$ and m machines.

Let us denote by T^* the optimal makespan for this instance.

Let us consider a list schedule of makespan T . Let us denote by $q(t)$ the number of active processors at time t .

We have $\forall t_1, t_2 \in [0, T] : t_1 \leq t_2 - T^* \Rightarrow q(t_1) + q(t_2) > m$ (otherwise, the tasks running at time t_2 could have been run at time t_1).

List Scheduling for Parallel Rigid Tasks

Let us assume we have n independent rigid jobs $J_1 = (p_1, q_1), \dots, J_n = (p_n, q_n)$ and m machines.

Let us denote by T^* the optimal makespan for this instance.

Let us consider a list schedule of makespan T . Let us denote by $q(t)$ the number of active processors at time t .

We have $\forall t_1, t_2 \in [0, T] : t_1 \leq t_2 - T^* \Rightarrow q(t_1) + q(t_2) > m$ (otherwise, the tasks running at time t_2 could have been run at time t_1).

Let us assume that $T > 2T^*$. Then we have:

$$\begin{aligned} mT^* &\geq \sum_i q_i p_i = \int_0^T q(t) = \int_0^{2T^*} q(t) + \int_{2T^*}^T q(t) \\ &\geq \underbrace{\int_0^{T^*} q(t) + q(t + T^*)}_{> mT^*} + \underbrace{\int_{2T^*}^T q(t)}_{\geq 0}, \text{ which is absurd.} \end{aligned}$$

List Scheduling for Parallel Rigid Tasks

Let us assume we have n independent rigid jobs $J_1 = (p_1, q_1), \dots, J_n = (p_n, q_n)$ and m machines.

Let us denote by T^* the optimal makespan for this instance.

Let us consider a list schedule of makespan T . Let us denote by $q(t)$ the number of active processors at time t .

We have $\forall t_1, t_2 \in [0, T] : t_1 \leq t_2 - T^* \Rightarrow q(t_1) + q(t_2) > m$ (otherwise, the tasks running at time t_2 could have been run at time t_1).

Let us assume that $T > 2T^*$. Then we have:

$$mT^* \geq \sum_i q_i p_i = \int_0^T q(t) dt = \int_0^{2T^*} q(t) dt + \int_{2T^*}^T q(t) dt$$

Theorem 2.

List-scheduling has an approximation factor of 2 for minimizing the C_{\max} of Parallel Rigid Tasks.

How can we use the previous result when going online?

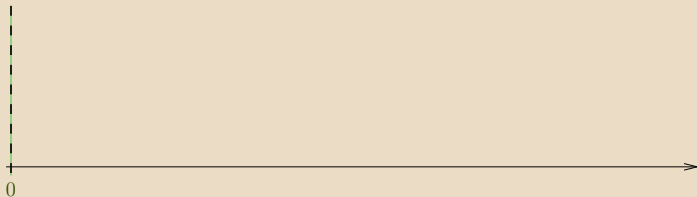
Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P | size_j | C_{\max} \rangle$.
Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P | size_j, r_j | C_{\max} \rangle$.

Proof.

Let us look at the schedule produced by \mathcal{A} on an instance \mathcal{I} .

release of
 S_0 jobs



How can we use the previous result when going online?

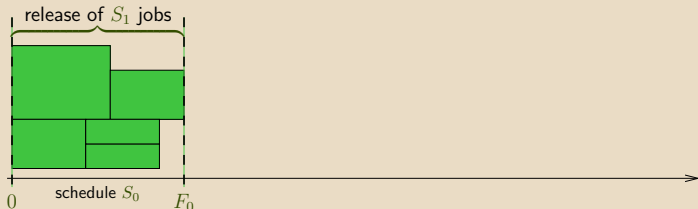
Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P | size_j | C_{\max} \rangle$.
Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P | size_j, r_j | C_{\max} \rangle$.

Proof.

Let us look at the schedule produced by \mathcal{A} on an instance \mathcal{I} .

release of
 S_0 jobs



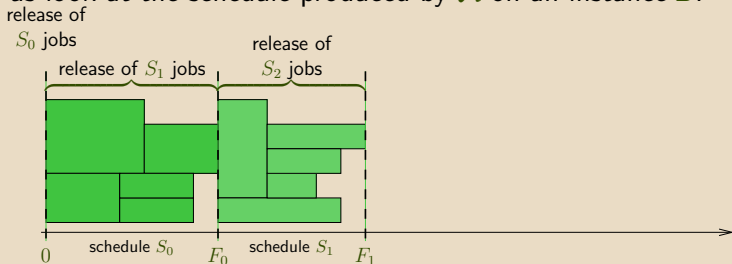
How can we use the previous result when going online?

Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P | size_j | C_{\max} \rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P | size_j, r_j | C_{\max} \rangle$.

Proof.

Let us look at the schedule produced by \mathcal{A} on an instance \mathcal{I} .



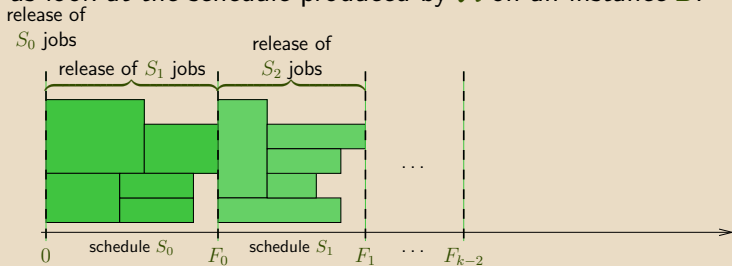
How can we use the previous result when going online?

Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P | size_j | C_{\max} \rangle$.
Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P | size_j, r_j | C_{\max} \rangle$.

Proof.

Let us look at the schedule produced by \mathcal{A} on an instance \mathcal{I} .



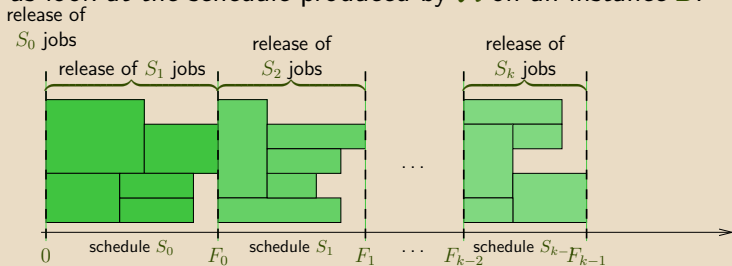
How can we use the previous result when going online?

Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P | size_j | C_{\max} \rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P | size_j, r_j | C_{\max} \rangle$.

Proof.

Let us look at the schedule produced by \mathcal{A} on an instance \mathcal{I} .



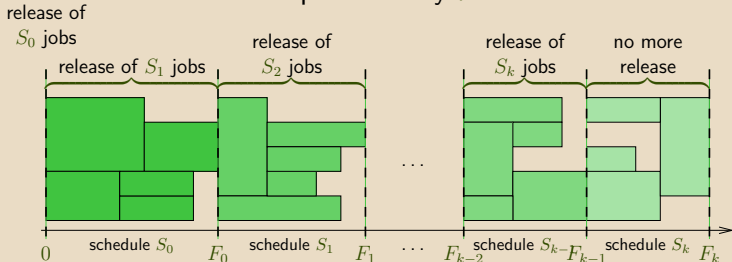
How can we use the previous result when going online?

Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P | size_j | C_{\max} \rangle$.
Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P | size_j, r_j | C_{\max} \rangle$.

Proof.

Let us look at the schedule produced by \mathcal{A} on an instance \mathcal{I} .



How can we use the previous result when going online?

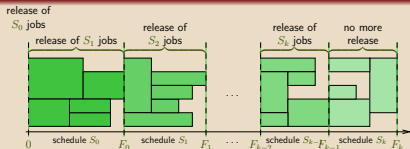
Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P|size_j|C_{\max}\rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P|size_j, r_j|C_{\max}\rangle$.

Proof.

Consider \mathcal{I}' where S_k jobs are released at time F_{k-2} . We have:

$$C_{\max}^*(\mathcal{I}') \leq C_{\max}^*(\mathcal{I}).$$



How can we use the previous result when going online?

Theorem 3: [Shmoys91].

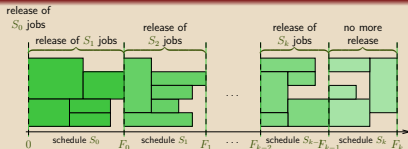
Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P|size_j|C_{\max}\rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P|size_j, r_j|C_{\max}\rangle$.

Proof.

Consider \mathcal{I}' where S_k jobs are released at time F_{k-2} . We have:

$$C_{\max}^*(\mathcal{I}') \leq C_{\max}^*(\mathcal{I}).$$

$$\blacktriangleright F_{k-2} + F_k - F_{k-1} \leq \rho C_{\max}^*(\mathcal{I}')$$



How can we use the previous result when going online?

Theorem 3: [Shmoys91].

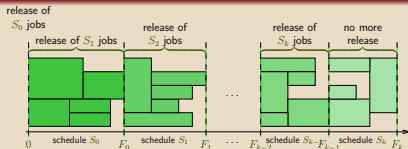
Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P|size_j|C_{\max}\rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P|size_j, r_j|C_{\max}\rangle$.

Proof.

Consider \mathcal{I}' where S_k jobs are released at time F_{k-2} . We have:

$$C_{\max}^*(\mathcal{I}') \leq C_{\max}^*(\mathcal{I}).$$

- ▶ $F_{k-2} + F_k - F_{k-1} \leq \rho C_{\max}^*(\mathcal{I}')$
- ▶ $F_{k-1} - F_{k-2} \leq \rho C_{\max}^*(\mathcal{I}')$



How can we use the previous result when going online?

Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P|size_j|C_{\max}\rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P|size_j, r_j|C_{\max}\rangle$.

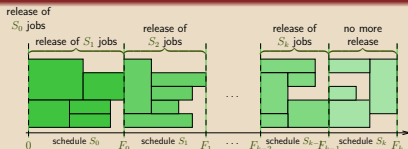
Proof.

Consider \mathcal{I}' where S_k jobs are released at time F_{k-2} . We have:

$$C_{\max}^*(\mathcal{I}') \leq C_{\max}^*(\mathcal{I}).$$

- ▶ $F_{k-2} + F_k - F_{k-1} \leq \rho C_{\max}^*(\mathcal{I}')$
- ▶ $F_{k-1} - F_{k-2} \leq \rho C_{\max}^*(\mathcal{I}')$

Hence $F_k \leq 2\rho C_{\max}^*(\mathcal{I}') \leq 2\rho C_{\max}^*(\mathcal{I})$



How can we use the previous result when going online?

Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P|size_j|C_{\max}\rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P|size_j, r_j|C_{\max}\rangle$.

- ▶ There is a PTAS for $\langle Q||C_{\max}\rangle$. Hence, there is an $(2 + \varepsilon)$ -competitive online clairvoyant algorithm for $\langle Q|r_j|C_{\max}\rangle$.

How can we use the previous result when going online?

Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P|size_j|C_{\max}\rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P|size_j, r_j|C_{\max}\rangle$.

- ▶ There is a PTAS for $\langle Q||C_{\max}\rangle$. Hence, there is an $(2 + \varepsilon)$ -competitive online clairvoyant algorithm for $\langle Q|r_j|C_{\max}\rangle$.
- ▶ There is a 2 approximation $\langle Q||C_{\max}\rangle$. Hence, there is an 4-competitive online clairvoyant algorithm for $\langle Q|r_j|C_{\max}\rangle$.

How can we use the previous result when going online?

Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P|size_j|C_{\max}\rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P|size_j, r_j|C_{\max}\rangle$.

- ▶ There is a PTAS for $\langle Q||C_{\max}\rangle$. Hence, there is an $(2 + \varepsilon)$ -competitive online clairvoyant algorithm for $\langle Q|r_j|C_{\max}\rangle$.
- ▶ There is a 2 approximation $\langle Q||C_{\max}\rangle$. Hence, there is an 4-competitive online clairvoyant algorithm for $\langle Q|r_j|C_{\max}\rangle$.
- ▶ There is a 2 approximation $\langle P|size_j|C_{\max}\rangle$. Hence, there is an 4-competitive online clairvoyant algorithm for $\langle Q|size_j|C_{\max}\rangle$.

How can we use the previous result when going online?

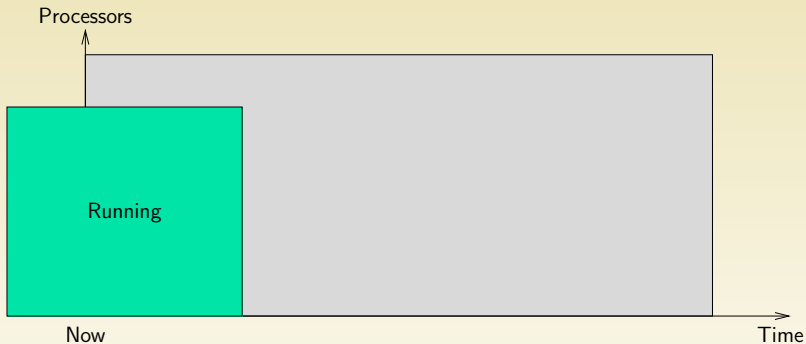
Theorem 3: [Shmoys91].

Let \mathcal{A} be a polynomial-time ρ -approximation for $\langle P|size_j|C_{\max}\rangle$. Based on \mathcal{A} , we can build a 2ρ -competitive polynomial-time online clairvoyant algorithm for $\langle P|size_j, r_j|C_{\max}\rangle$.

- ▶ There is a PTAS for $\langle Q||C_{\max}\rangle$. Hence, there is an $(2 + \varepsilon)$ -competitive online clairvoyant algorithm for $\langle Q|r_j|C_{\max}\rangle$.
- ▶ There is a 2 approximation $\langle Q||C_{\max}\rangle$. Hence, there is an 4-competitive online clairvoyant algorithm for $\langle Q|r_j|C_{\max}\rangle$.
- ▶ There is a 2 approximation $\langle P|size_j|C_{\max}\rangle$. Hence, there is an 4-competitive online clairvoyant algorithm for $\langle Q|size_j|C_{\max}\rangle$.
- ▶ Actually, by doing a slightly finer analysis, one can show that the list-scheduling algorithm is a $(2 - 1/m)$ -competitive non-clairvoyant algorithm for $\langle P|r_j|C_{\max}\rangle$.

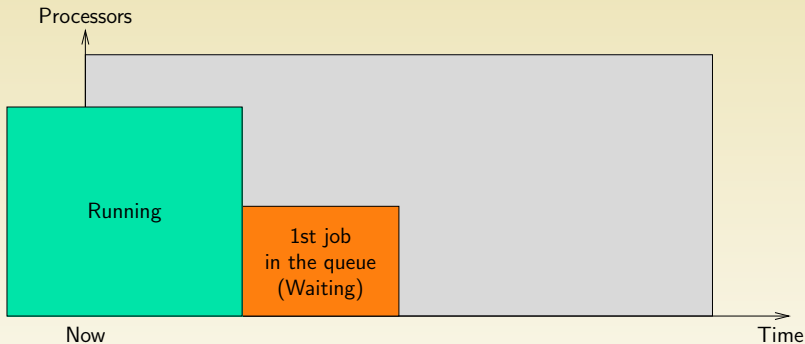
- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 **Batch Scheduling**
 - Principles
 - Theoretical results
 - **Basic idea: FCFS + Backfilling**
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

General Principle



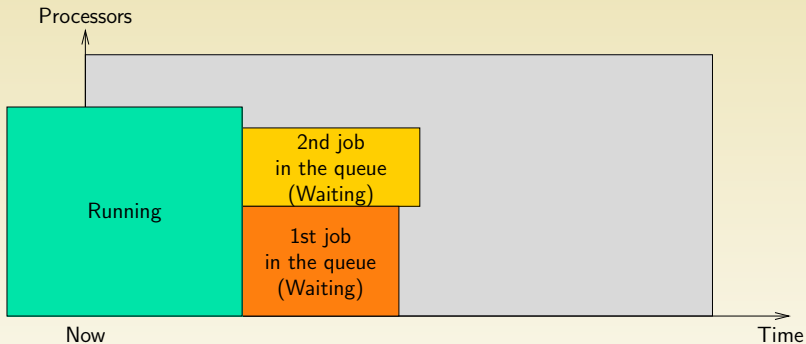
- ▶ Jobs arrive one after the other and are scheduled at arrival.

General Principle



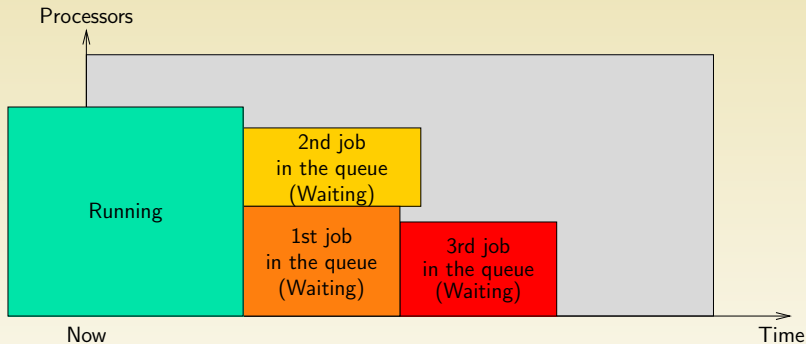
- ▶ Jobs arrive one after the other and are scheduled at arrival.

General Principle



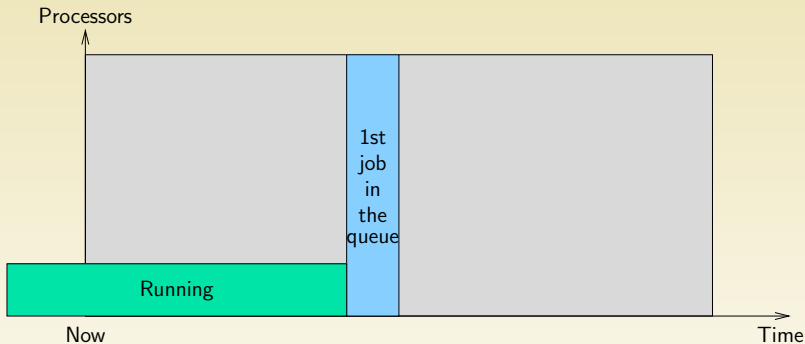
- ▶ Jobs arrive one after the other and are scheduled at arrival.

General Principle



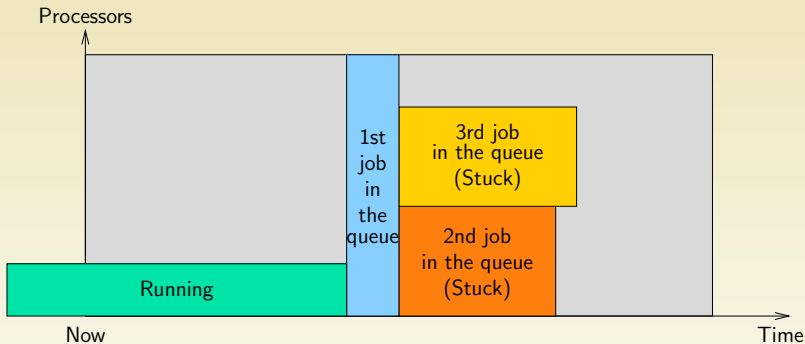
- ▶ Jobs arrive one after the other and are scheduled at arrival.

First Come First Served



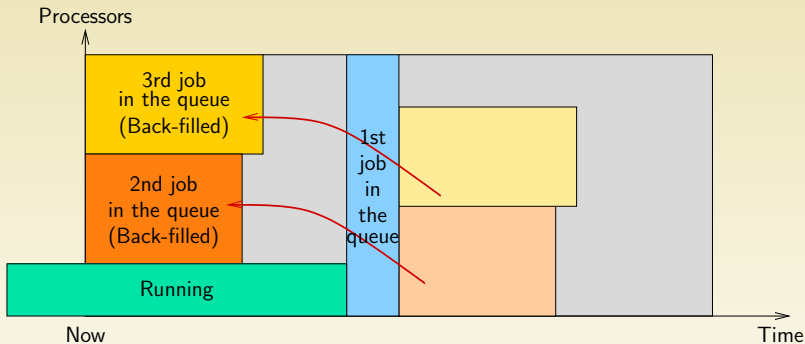
- ▶ FCFS = simplest scheduling option
- ▶ Fragmentation \rightsquigarrow need for **backfilling**

First Come First Served



- ▶ FCFS = simplest scheduling option
- ▶ Fragmentation \leadsto need for **backfilling**

First Come First Served



- ▶ FCFS = simplest scheduling option
- ▶ Fragmentation \rightsquigarrow need for **backfilling**

- ▶ Which job(s) should be picked for promotion through the queue?
- ▶ Many heuristics are possible
- ▶ Two have been studied in detail
 - ▶ EASY
 - ▶ Conservative Back Filling (CBF)
- ▶ In practice EASY (or variants of it) is used, while CBF is not.
- ▶ Although, OAR, a recently proposed batch scheduler implements CBF.

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - **EASY**
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Extensible Argonne Scheduling System

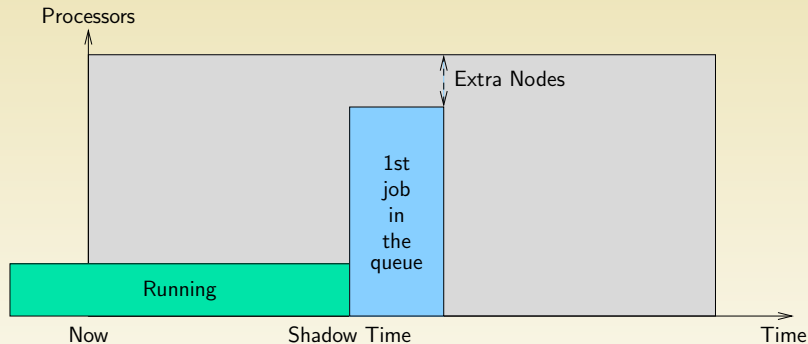
Maintain only one *reservation*, for the first job in the queue.

Definitions:

Shadow time time at which the first job in the queue starts execution

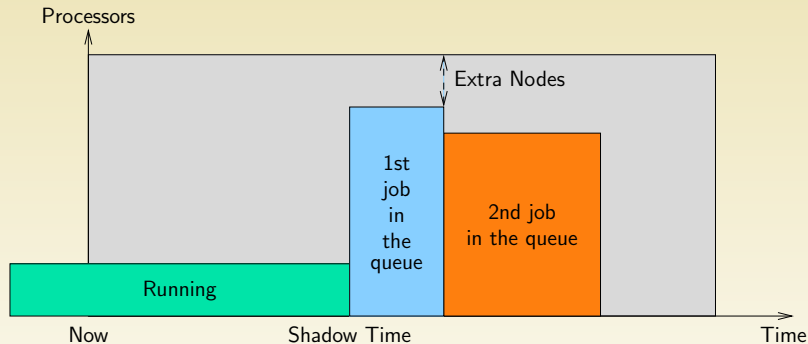
Extra nodes number of nodes idle when the first job in the queue starts execution

- 1 Go through the queue in order starting with the 2nd job.
- 2 Backfill a job if it will terminate by the shadow time, **or** it needs less than the extra nodes.



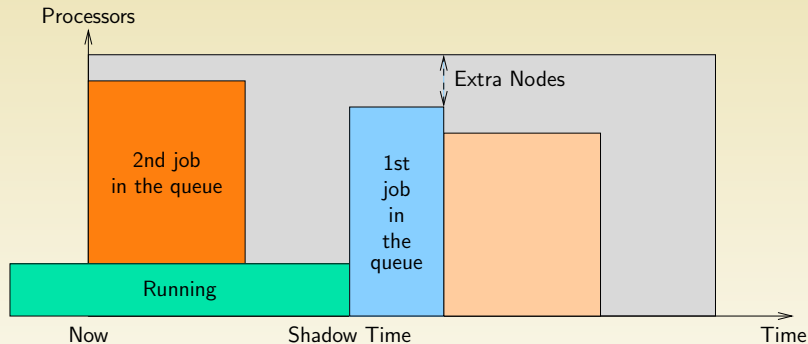
Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs



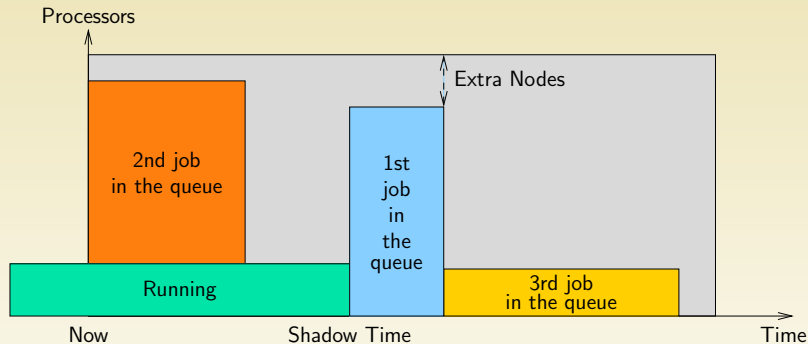
Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs



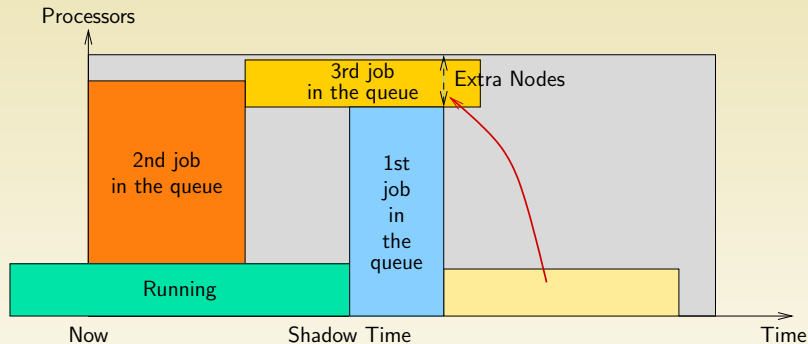
Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs



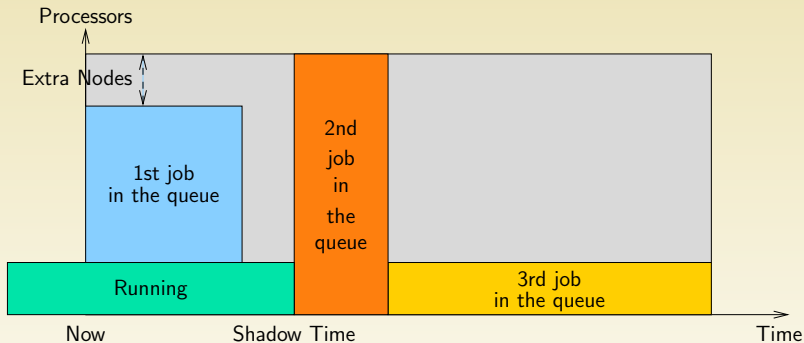
Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs



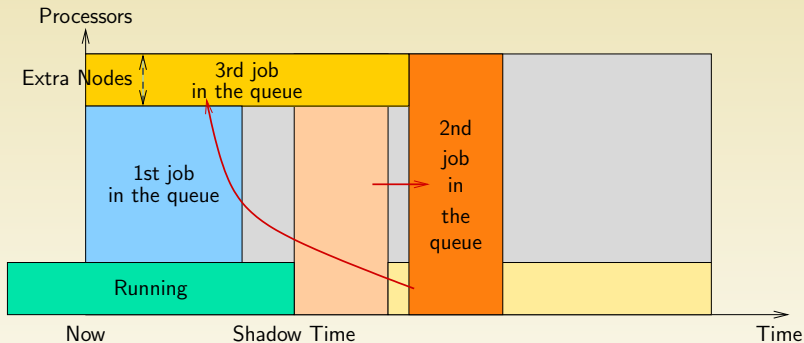
Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs



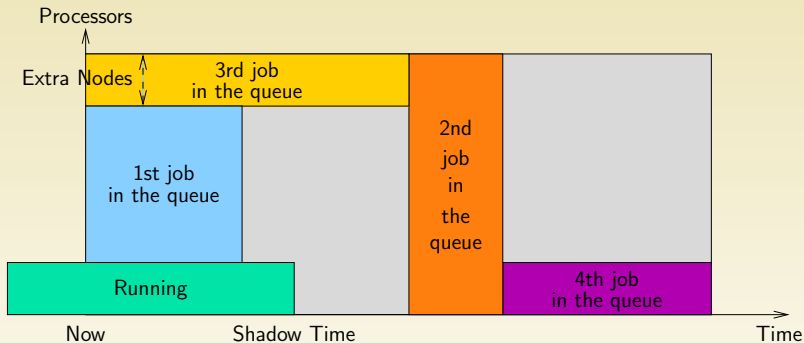
Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs
- ▶ BUT, other jobs may be delayed infinitely!



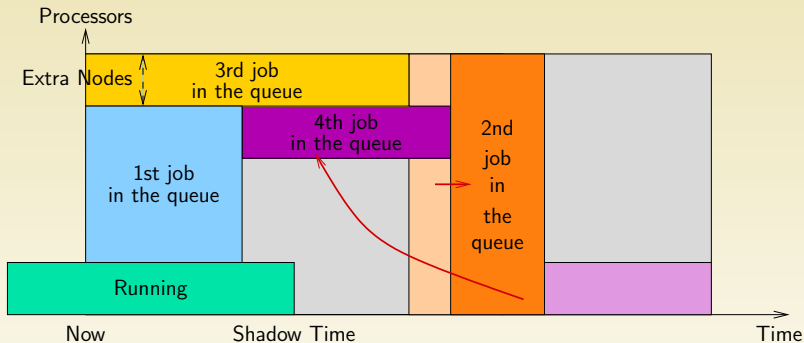
Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs
- ▶ BUT, other jobs may be delayed infinitely!



Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs
- ▶ BUT, other jobs may be delayed infinitely!



Property:

- ▶ The first job in the queue will never be delayed by backfilled jobs
- ▶ BUT, other jobs may be delayed infinitely!

Unbounded Delay. ▶ The first job in the queue will never be delayed by backfilled jobs

- ▶ BUT, other jobs may be delayed infinitely!

No Starvation. ▶ Delay of first job is bounded by runtime of current jobs

- ▶ When the first job finishes, the second job becomes the first job in the queue
- ▶ Once it is the first job, it cannot be delayed further

Other approach. ▶ **Conservative Backfilling.** *EVERY* job has a *reservation*. A job may be backfilled only if it does not delay any other job ahead of it in the queue.

- ▶ Fixes the unbounded delay problem that EASY has. More complicated to implement (The algorithm must find holes in the schedule) though.
- ▶ EASY favors small long jobs and harms large short jobs.

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 **Batch Scheduling**
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - **How Good is the Schedule?**
- 4 Gang Scheduling as an Alternative
 - Principles

When Does Backfilling Happen?

Possibly when

- ▶ A new job arrives

When Does Backfilling Happen?

Possibly when

- ▶ A new job arrives
- ▶ The first job in the queue starts

When Does Backfilling Happen?

Possibly when

- ▶ A new job arrives
- ▶ The first job in the queue starts
- ▶ **When a job finishes early**

When Does Backfilling Happen?

Possibly when

- ▶ A new job arrives
- ▶ The first job in the queue starts
- ▶ **When a job finishes early**

Users provide job runtime **estimates** (Jobs are killed if they go over).

When Does Backfilling Happen?

Possibly when

- ▶ A new job arrives
- ▶ The first job in the queue starts
- ▶ **When a job finishes early**

Users provide job runtime **estimates** (Jobs are killed if they go over).

Trade-off:

- ▶ provide a **conservative estimate**: you goes through the queue faster (may be backfilled)
- ▶ provide a **loose estimate**: your job will not be killed

When Does Backfilling Happen?

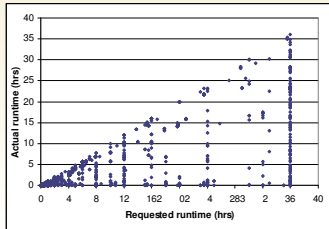
Possibly when

- ▶ A new job arrives
- ▶ The first job in the queue starts
- ▶ **When a job finishes early**

Users provide job runtime **estimates** (Jobs are killed if they go over).
Trade-off:

- ▶ provide a **conservative estimate**: you goes through the queue faster (may be backfilled)
- ▶ provide a **loose estimate**: your job will not be killed

Are estimates accurate?



How Good is the Schedule ?

All of this is great, but how do we know what a “good” schedule is? FCFS, EASY, CFB, Random?

What we need are **metrics** to quantify how good a schedule is. It has to be an aggregate metric over all jobs

How Good is the Schedule ?

All of this is great, but how do we know what a “good” schedule is? FCFS, EASY, CFB, Random?

What we need are **metrics** to quantify how good a schedule is. It has to be an aggregate metric over all jobs

- 1 Turn-around time or **flow** (Wait time + Run time).

Job 1 needs 1h of compute time and waits 1s

Job 2 needs 1s of compute time and waits 1h

Clearly Job 1 is really happy, and Job 2 is not happy at all

How Good is the Schedule ?

All of this is great, but how do we know what a “good” schedule is? FCFS, EASY, CFB, Random?

What we need are **metrics** to quantify how good a schedule is. It has to be an aggregate metric over all jobs

- 1 Turn-around time or **flow** (Wait time + Run time).

Job 1 needs 1h of compute time and waits 1s

Job 2 needs 1s of compute time and waits 1h

Clearly Job 1 is really happy, and Job 2 is not happy at all

- 2 **Wait time** (equivalent to “user happiness”)

Job 1 asks for 1 nodes and waits 1 h

Job 2 asks for 512 nodes and waits 1h

Again, Job 1 is unhappy while Job 2 is probably sort of happy.

We need a metric that represents happiness for small, large, short, long jobs.

How Good is the Schedule ?

All of this is great, but how do we know what a “good” schedule is? FCFS, EASY, CFB, Random?

What we need are **metrics** to quantify how good a schedule is. It has to be an aggregate metric over all jobs

- 1 Turn-around time or **flow** (Wait time + Run time).

Job 1 needs 1h of compute time and waits 1s

Job 2 needs 1s of compute time and waits 1h

Clearly Job 1 is really happy, and Job 2 is not happy at all

- 2 **Wait time** (equivalent to “user happiness”)

Job 1 asks for 1 nodes and waits 1 h

Job 2 asks for 512 nodes and waits 1h

Again, Job 1 is unhappy while Job 2 is probably sort of happy.

We need a metric that represents happiness for small, large, short, long jobs.

- 3 Slowdown or **Stretch** (turn-around time divided by turn-around time if alone in the system)

Doesn't really take care of the small/large problem. Could think of some scaling, but unclear !

Now we have a few metrics we can consider

We can run simulations of the scheduling algorithms, and see how they fare.

We need to test these algorithms in representative scenarios

Supercomputer/cluster traces. Collect the following for long periods of time:

- ▶ Time of submission
- ▶ How many nodes asked
- ▶ How much time asked
- ▶ How much time was actually used
- ▶ How much time spent in the queue

Uses of the traces:

- 1 Drive simulations
- 2 Come up with models of user behaviors

Sample Results

A type of experiments that people have done: replace user estimate by f times the actual run time

Possible to improve performance by multiplying user estimates by 2!

	EASY	CBF
Mean Slowdown		
KTH	-4.8%	-23.0%
CTC	-7.9%	-18.0%
SDSC	+4.6%	-14.2%
Mean Response time		
KTH	-3.3%	-7.0%
CTC	-0.9%	-1.6%
SDSC	-1.6%	-10.9%

- ▶ These are all **heuristics**.
- ▶ They are not specifically designed to optimize the metrics we have designed.
- ▶ It is difficult to truly understand the reasons for the results.
- ▶ But one can derive some empirical wisdom.
- ▶ One of the reasons why one is stuck with possibly obscure heuristics is that we're dealing with an *on-line* problem: We don't know what happens next.
- ▶ We cannot wait for all jobs to be submitted to make a decision. But we can wait for a while, accumulate jobs, and schedule them together.

Batch Schedulers are what we're stuck with at the moment.
They are often hated by users.

- ▶ I submit to the queue asking for 10 nodes for 1 hour.
- ▶ I wait for two days.
- ▶ My code finally starts, but doesn't finish within 1 hour and gets killed!!

Batch Schedulers are what we're stuck with at the moment. They are often hated by users.

- ▶ I submit to the queue asking for 10 nodes for 1 hour.
- ▶ I wait for two days.
- ▶ My code finally starts, but doesn't finish within 1 hour and gets killed!!

A lot of research, a few things happening “in the field”. When you go to a company that has clusters (like most of them), they typically have a job scheduler, so it's good to have some idea of what it is.

Batch Schedulers are what we're stuck with at the moment. They are often hated by users.

- ▶ I submit to the queue asking for 10 nodes for 1 hour.
- ▶ I wait for two days.
- ▶ My code finally starts, but doesn't finish within 1 hour and gets killed!!

A lot of research, a few things happening “in the field”.

When you go to a company that has clusters (like most of them), they typically have a job scheduler, so it's good to have some idea of what it is.

A completely different approach is **gang scheduling**, which we discuss next.

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

- ▶ All processes belonging to a job run at the same time (the term **gang** denotes all processors within a job).
- ▶ Each process runs alone on each processor.
- ▶ BUT: there is rapid **coordinated** context switching.
- ▶ It is possible to **suspend/preempt** jobs arbitrarily

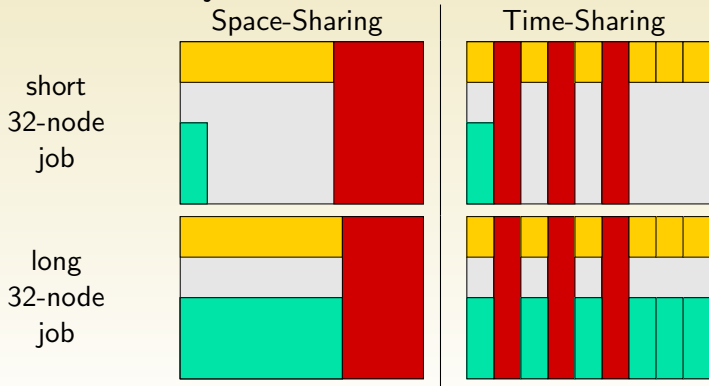
- ▶ All processes belonging to a job run at the same time (the term **gang** denotes all processors within a job).
- ▶ Each process runs alone on each processor.
- ▶ BUT: there is rapid **coordinated** context switching.
- ▶ It is possible to **suspend/preempt** jobs arbitrarily \rightsquigarrow May allow more flexibility to optimize some metrics.

- ▶ All processes belonging to a job run at the same time (the term **gang** denotes all processors within a job).
- ▶ Each process runs alone on each processor.
- ▶ BUT: there is rapid **coordinated** context switching.
- ▶ It is possible to **suspend/preempt** jobs arbitrarily \rightsquigarrow May allow more flexibility to optimize some metrics.
- ▶ If processing times are not known in advance (or grossly erroneous), preemption can help short jobs that would be “stuck” behind a long job.
- ▶ Should improve machine utilization.

Gang Scheduling: an Example

- ▶ A 128 node cluster.
- ▶ A running 64-node job.
- ▶ A 32-node job and a 128-node job are queued.

Should the 32-node job be started ?



More uniform slowdown, better resource usage.

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

- ▶ Overhead for context switching (trade-off between overhead and fine grain).

Gang Scheduling: Drawbacks

- ▶ Overhead for context switching (trade-off between overhead and fine grain).
- ▶ Overhead for coordinating context switching across multiple processors.

Gang Scheduling: Drawbacks

- ▶ Overhead for context switching (trade-off between overhead and fine grain).
- ▶ Overhead for coordinating context switching across multiple processors.
- ▶ Reduced cache efficiency(Frequent cache flushing).

Gang Scheduling: Drawbacks

- ▶ Overhead for context switching (trade-off between overhead and fine grain).
- ▶ Overhead for coordinating context switching across multiple processors.
- ▶ Reduced cache efficiency(Frequent cache flushing).
- ▶ RAM Pressure (more jobs must fit in memory, swapping to disk causes unacceptable overhead).

Gang Scheduling: Drawbacks

- ▶ Overhead for context switching (trade-off between overhead and fine grain).
- ▶ Overhead for coordinating context switching across multiple processors.
- ▶ Reduced cache efficiency(Frequent cache flushing).
- ▶ RAM Pressure (more jobs must fit in memory, swapping to disk causes unacceptable overhead).
- ▶ Typically not used in production HPC systems (batch scheduling is preferred).

Gang Scheduling: Drawbacks

- ▶ Overhead for context switching (trade-off between overhead and fine grain).
- ▶ Overhead for coordinating context switching across multiple processors.
- ▶ Reduced cache efficiency(Frequent cache flushing).
- ▶ RAM Pressure (more jobs must fit in memory, swapping to disk causes unacceptable overhead).
- ▶ Typically not used in production HPC systems (batch scheduling is preferred).
- ▶ Some implementations (MOSIX, Kerighed).

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Batch Scheduling it is then

So it seems we're stuck with batch scheduling.
Why don't we like Batch Scheduling?

So it seems we're stuck with batch scheduling.

Why don't we like Batch Scheduling? Because queue waiting times are difficult to predict.

- ▶ depends on the status of the queue
- ▶ depends on the scheduling algorithm used
- ▶ depends on all sorts of configuration parameters set by system administrator
- ▶ depends on future job completions!
- ▶ etc.

So I submit my job and then it's in limbo somewhere, which is eminently annoying to most users.

Batch Scheduling it is then

So it seems we're stuck with batch scheduling.

Why don't we like Batch Scheduling? Because queue waiting times are difficult to predict.

- ▶ depends on the status of the queue
- ▶ depends on the scheduling algorithm used
- ▶ depends on all sorts of configuration parameters set by system administrator
- ▶ depends on future job completions!
- ▶ etc.

So I submit my job and then it's in limbo somewhere, which is eminently annoying to most users.

That is why there is more and more demand for **reservation** support. Users build (badly?) the schedule by themselves.

- 1 Modeling Applications, General Notions
 - Introducing Fundamental Notions Through the Matrix Product Example
 - Adaptive Parallel Programs
 - Task Graphs and Parallel Tasks From Outer Space
- 2 Defining a Scheduling Problem
 - Rules of the Game
 - Criteria: How Do You Win the Game?
 - Analysis Method
 - Graham Notation
- 3 Batch Scheduling
 - Principles
 - Theoretical results
 - Basic idea: FCFS + Backfilling
 - EASY
 - How Good is the Schedule?
- 4 Gang Scheduling as an Alternative
 - Principles

Batch Scheduling and Grids

Grids result from the **collaboration** of many Universities/Computing Centers.

Everyone runs its own Batch Scheduler that cannot be bypassed.
How to decide where we should submit our jobs?

Batch Scheduling and Grids

Grids result from the **collaboration** of many Universities/Computing Centers.

Everyone runs its own Batch Scheduler that cannot be bypassed.
How to decide where we should submit our jobs?

When in doubt, a brute-force approach is to:

- ▶ Do multiple submissions for different numbers of nodes
- ▶ Cancel all submissions but the first one that comes back
- ▶ Or possibly make some ad-hoc call regarding whether to keep a potentially poor request in the hope of getting a better one through shortly after.

Grids result from the **collaboration** of many Universities/Computing Centers.

Everyone runs its own Batch Scheduler that cannot be bypassed.
How to decide where we should submit our jobs?

When in doubt, a brute-force approach is to:

- ▶ Do multiple submissions for different numbers of nodes
- ▶ Cancel all submissions but the first one that comes back
- ▶ Or possibly make some ad-hoc call regarding whether to keep a potentially poor request in the hope of getting a better one through shortly after.

What happens if everybody does this?

Grids result from the **collaboration** of many Universities/Computing Centers.

Everyone runs its own Batch Scheduler that cannot be bypassed.
How to decide where we should submit our jobs?

When in doubt, a brute-force approach is to:

- ▶ Do multiple submissions for different numbers of nodes
- ▶ Cancel all submissions but the first one that comes back
- ▶ Or possibly make some ad-hoc call regarding whether to keep a potentially poor request in the hope of getting a better one through shortly after.

What happens if everybody does this?

Other issues:

- ▶ File Staging ?
- ▶ Load Balancing between sites ?

Sequential Job Scheduling for Grids

A set unrelated processors P_1, \dots, P_n and a set of sequential jobs J_1, \dots, J_n (processing time $p_{i,j}$).

Let's try a few natural scheduling strategies. We denote by a_i the time at which P_i is available (at the beginning $a_i = 0$ for all P_i):

Sequential Job Scheduling for Grids

A set unrelated processors P_1, \dots, P_n and a set of sequential jobs J_1, \dots, J_n (processing time $p_{i,j}$).

Let's try a few natural scheduling strategies. We denote by a_i the time at which P_i is available (at the beginning $a_i = 0$ for all P_i):

Min-Min Compute the minimum completion time $C_j = a_i + p_{i,j}$ of each J_j and choose the one with the smallest C_j . Update the corresponding a_i (its best host) accordingly ($a_i \leftarrow a_i + p_{i,j}$).

Sequential Job Scheduling for Grids

A set unrelated processors P_1, \dots, P_n and a set of sequential jobs J_1, \dots, J_n (processing time $p_{i,j}$).

Let's try a few natural scheduling strategies. We denote by a_i the time at which P_i is available (at the beginning $a_i = 0$ for all P_i):

Min-Min Compute the minimum completion time $C_j = a_i + p_{i,j}$ of each J_j and choose the one with the smallest C_j . Update the corresponding a_i (its best host) accordingly ($a_i \leftarrow a_i + p_{i,j}$).

Max-Min Choose J_j with the largest C_j and update the corresponding a_i (its best host) accordingly.

Sequential Job Scheduling for Grids

A set unrelated processors P_1, \dots, P_n and a set of sequential jobs J_1, \dots, J_n (processing time $p_{i,j}$).

Let's try a few natural scheduling strategies. We denote by a_i the time at which P_i is available (at the beginning $a_i = 0$ for all P_i):

Min-Min Compute the minimum completion time $C_j = a_i + p_{i,j}$ of each J_j and choose the one with the smallest C_j . Update the corresponding a_i (its best host) accordingly ($a_i \leftarrow a_i + p_{i,j}$).

Max-Min Choose J_j with the largest C_j and update the corresponding a_i (its best host) accordingly.

Sufferage S_j is the difference between the best completion time of J_j and its second best completion time. Choose the job with the largest sufferage and schedule it on its best processor.

Sequential Job Scheduling for Grids

A set unrelated processors P_1, \dots, P_n and a set of sequential jobs J_1, \dots, J_n (processing time $p_{i,j}$).

Let's try a few natural scheduling strategies. We denote by a_i the time at which P_i is available (at the beginning $a_i = 0$ for all P_i):

Min-Min Compute the minimum completion time $C_j = a_i + p_{i,j}$ of each J_j and choose the one with the smallest C_j . Update the corresponding a_i (its best host) accordingly ($a_i \leftarrow a_i + p_{i,j}$).

Max-Min Choose J_j with the largest C_j and update the corresponding a_i (its best host) accordingly.

Sufferage S_j is the difference between the best completion time of J_j and its second best completion time. Choose the job with the largest sufferage and schedule it on its best processor.

Problem: How do you get an estimate of $p_{i,j}$?

So Where are we ?

- ▶ Batch schedulers are complex pieces of software that are used in practice.
- ▶ A lot of experience on how they work and how to use them.
- ▶ But ultimately everybody knows they are an imperfect solution.
- ▶ Many view the lack of theoretical foundations as a big problem.
- ▶ Some just don't care. . .

Fools ignore complexity. Pragmatists suffer it. Some can avoid it. Geniuses remove it.

– *"Epigrams in Programming"*, by Alan J. Perlis of Yale University.

Bibliography



A.J. Bernstein.

Analysis of programs for parallel processing.

IEEE Transactions on Electronic Computers, 15:757–762, October 1966.



Peter Brucker.

Scheduling Algorithms.

Springer, Heidelberg, 2 edition, 1998.



E. G. Coffman.

Computer and job-shop scheduling theory.

John Wiley & Sons, 1976.



D.B. Shmoys, J. Wein, and D.P. Williamson.

Scheduling parallel machines on-line.

Symposium on Foundations of Computer Science, 0:131–140, 1991.