# Planning & Scheduling

**Roman Barták**

Department of Theoretical Computer Science and Mathematical Logic

**Constraint-based Scheduling**

---

## Scheduling model

- **Scheduling problem** is static so it can be directly **encoded as a CSP**.

- Constraint technology is used for **full scheduling**.

**CSP model:**

- **variables**
    - position of activity A in time and space
    - time allocation:    **start(A), [p(A), end(A)]**
    - resource allocation: **resource(A)**
- **Domain**
    - **release times** and **deadlines** for the time variables
    - **alternative resources** for the resource variables
- **constraints**
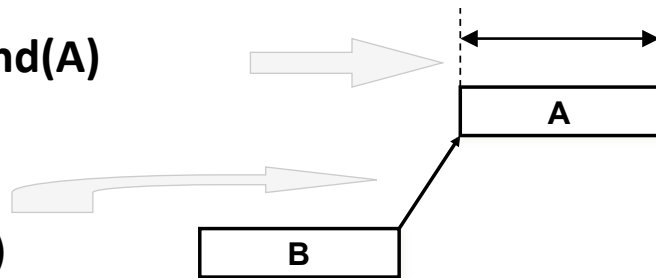    - sequencing and resource capacities

- **time relations**
  - **start(A)+p(A)=end(A)**
  - sequencing
    - B<<A
    - ↳ **end(B)≤start(A)**

↳ **resource capacity constraints**
  - unary resource (activities cannot overlap)
    - A<<B ∨ B<<A
    - ↳ **end(A)≤start(B) ∨ end(B)≤start(A)**

**Resources are used in slightly different meanings in planning and scheduling!**

- **scheduling**
  - resource
    = a **machine** (space) for processing the activity

- **planning**
  - resource
    = consumed/produced **material** by the activity
  - resource in the scheduling sense is often handled via logical precondition (e.g. hand is free)

- **Unary resource**
  - at most one activity can be processed at any time

- **Cumulative resource**
  - several activities can be processed in parallel if resource capacity is not exceeded.

- **Producible/consumable resource**
  - activity consumes/produces some quantity of the resource
  - minimal capacity is requested (consumption) and maximal capacity cannot be exceeded (production)

**Activities cannot overlap in time!**
  - at any time at most one activity can be processed, hence these resources are called **unary**
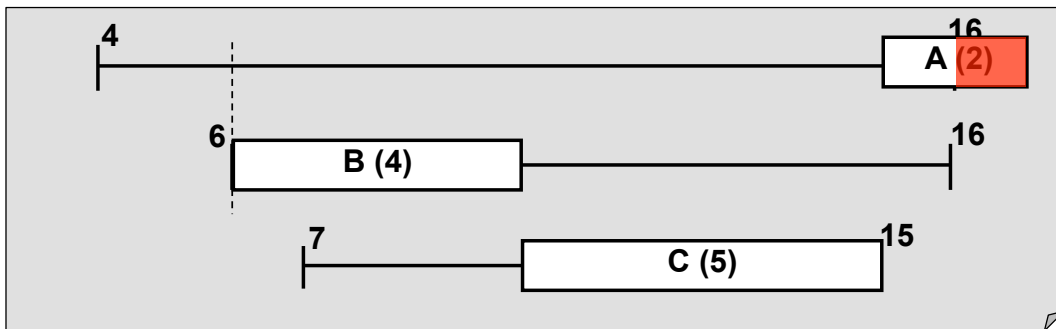
We assume that activities are **not interruptible** (non-preemptible).
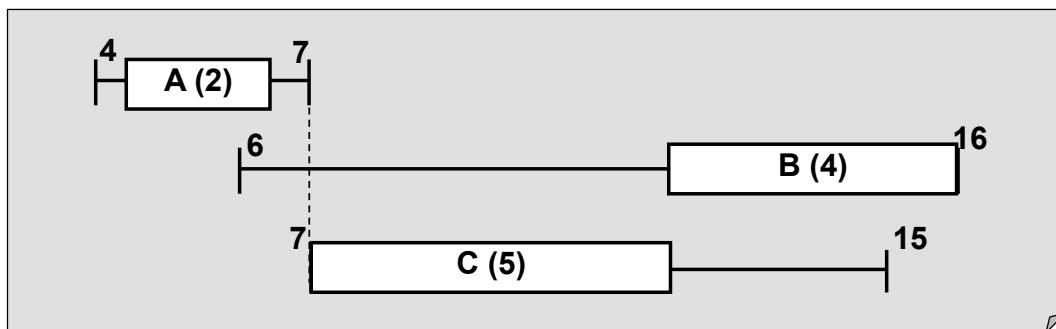  - non-interruptible (non-preemptible) activity occupies the resource from its start till its completion

A simple mode is based on disjunctive resources:
  - **A « B ∨ B « A**

    **end(A) ≤ start(B) ∨ end(B) ≤ start(A)**
  - hence these resources ale also called **disjunctive**

**What happens if activity A is not processed first?**



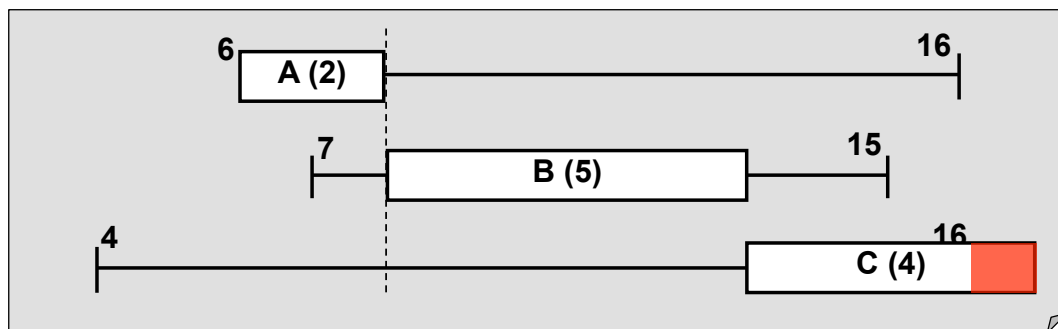**Not enough time for A, B, and C and thus A must be first!**

## Inference rules:

- **$p(\Omega \cup \{A\}) > lct(\Omega \cup \{A\}) - est(\Omega) \Rightarrow A«\Omega$**

- **$p(\Omega \cup \{A\}) > lct(\Omega) - est(\Omega \cup \{A\}) \Rightarrow \Omega«A$**

- **$A«\Omega \Rightarrow end(A) \leq min\{ lct(\Omega') - p(\Omega') \mid \Omega' \subseteq \Omega \}$**

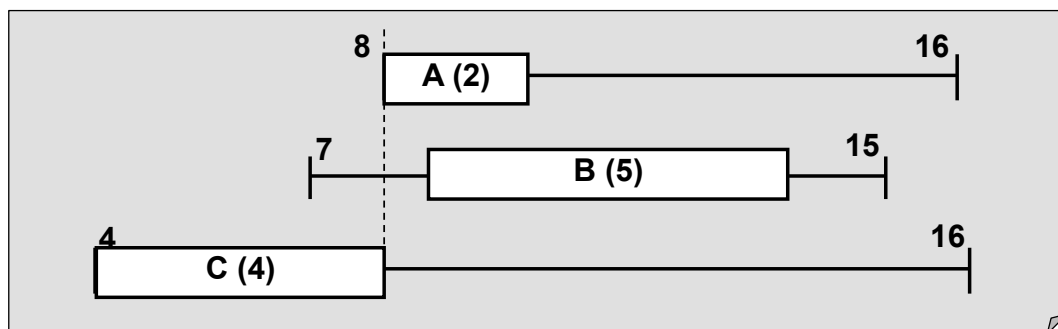- **$\Omega«A \Rightarrow start(A) \geq max\{ est(\Omega') + p(\Omega') \mid \Omega' \subseteq \Omega \}$**

## In practice:

- we need to explore $n.2^n$ pairs $(A,\Omega)$ (that is too many!)

- instead of all $\Omega$ we can use **task intervals** [X,Y]
  $\{C \mid est(X) \leq est(C) \wedge lct(C) \leq lct(Y)\}$

  ↳ time complexity is $O(n^3)$; this is frequently used incremental approach

  ↳ there are also algorithms with time complexity $O(n^2)$ and $O(n.\log n)$

**What happens if activity A is processed first?**



**Not enough time for B and C and thus A cannot be first!**

## Not-first inference rules:

$$\min(start(A)) + p(\Omega) + p(A) > \max(end(\Omega)) \Rightarrow \neg\, A{<}{<}\Omega$$

$$\neg\, A{<}{<}\Omega \Rightarrow start(A) \geq \min\{\, end(B) \mid B{\in}\Omega \,\}$$

## Not-last (symmetrical) inference rules:

$$\min(start(\Omega)) + p(\Omega) + p(A) > \max(end(A)) \Rightarrow \neg\, \Omega{<}{<}A$$

$$\neg\, \Omega{<}{<}A \Rightarrow end(A) \leq \max\{\, start(B) \mid B{\in}\Omega \,\}$$
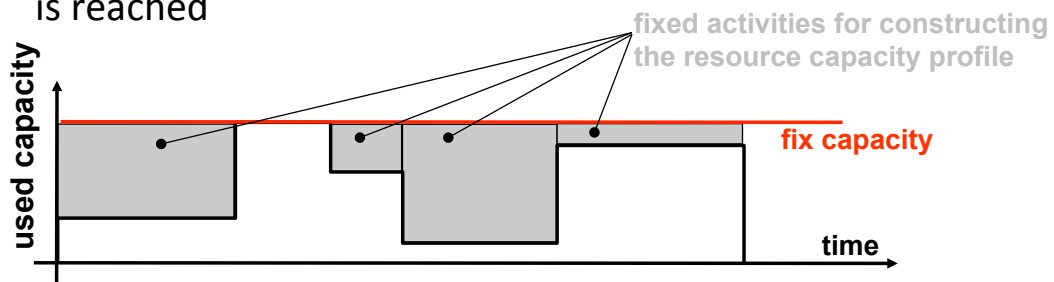
## In practice:

– it is possible to use selected sets $\Omega$ only

– time complexity $O(n^2)$

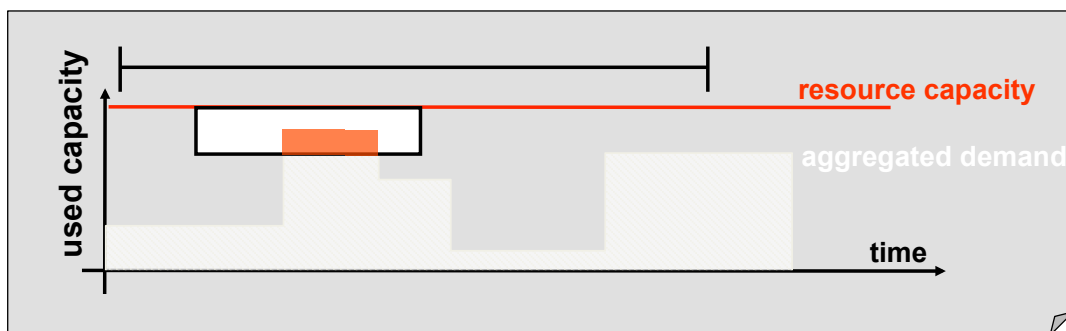Each **activity uses some capacity** of the resource – **cap(A)**.

Activities can be **processed in parallel** if a resource capacity is not exceeded.
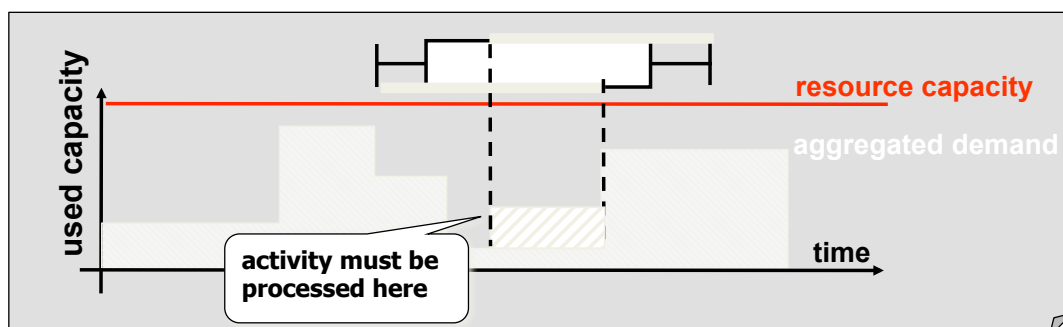
Resource capacity **may vary in time**

– modeled via fix capacity over time and fixed activities consuming the resource until the requested capacity level is reached



fixed activities for constructing the resource capacity profile

fix capacity

used capacity

time

**Where is enough capacity for processing the activity?**



resource capacity

aggregated demand

used capacity

time

**How the aggregated demand is constructed?**



resource capacity

aggregated demand

used capacity

activity must be processed here

time

- How to ensure that capacity is not exceed at any time point?*

$$\forall t \quad \sum_{start(A_i) \le t \le end(A_i)} cap(A_i) \le MaxCapacity$$

- **Timetable** for the activity A is a set of Boolean variables **X(A,t)** indicating whether A is processed in time t.
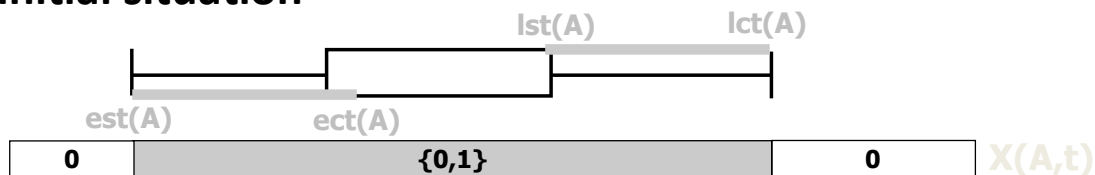
$$\forall t \quad \sum_{A_i} X(A_i,t) \cdot cap(A_i) \le MaxCapacity$$

$$\forall t,i \quad start(A_i) \le t \le end(A_i) \Leftrightarrow X(A_i,t)$$
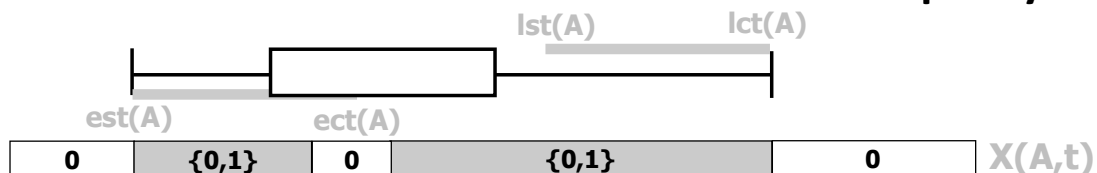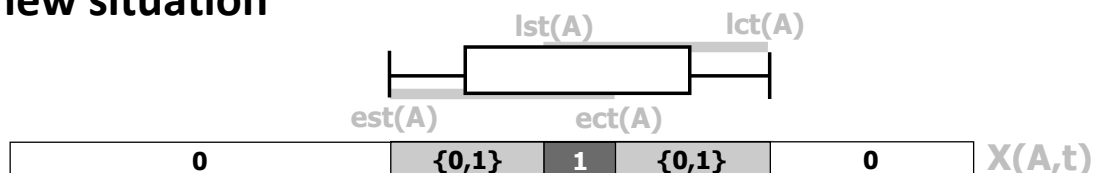
* discrete time is expected

**initial situation**



**some times are forbidden due to exceeded capacity**



**new situation**

**How to implement domain filtering?**

$$\forall t,i \ \ start(A_i) \leq t < end(A_i) \Leftrightarrow X(A_i,t)$$

Problem:

t is both an index and variable

**start(A) ≥ min{t : ub(X(A,t))=1}**

**end(A) ≤ 1+max{t : ub(X(A,t))=1}**

**X(A,t)=0 ∧ t<ect(A) ⟹ start(A)>t**

**X(A,t)=0 ∧ lst(A)≤t ⟹ end(A)≤t**

**(lst(A)≤t ∧ t<ect(A) ⟹ X(A,t)=1)**

- **How to model alternative resources for a given activity?**
- Use a **duplicate activity** for each resource.
  - duplicate activity participates in a respective resource constraint but does not restrict other activities there
    - „failure" means removing the resource from the domain of variable res(A)
    - deleting the resource from the domain of variable res(A) means „deleting" the respective duplicate activity
  - original activity participates in precedence constraints (e.g. within a job)
  - restricted times of duplicate activities are propagated to the original activity and vice versa.

Let $A_u$ be a duplicate of activity A allocated to resource $u \in$ res(A).

> **u∈res(A) $\Rightarrow$ start(A) ≤ start($A_u$)**
>
> **u∈res(A) $\Rightarrow$ end($A_u$) ≤ end(A)**
>
> **start(A) ≥ min{start($A_u$) : u∈res(A)}**
>
> **end(A) ≤ max{end($A_u$) : u∈ res(A)}**
>
> **empty time windows for $A_u \Rightarrow$ res(A)\\{u}**

> In practice, this inference is identical to constructive disjunction between the alternative resources.
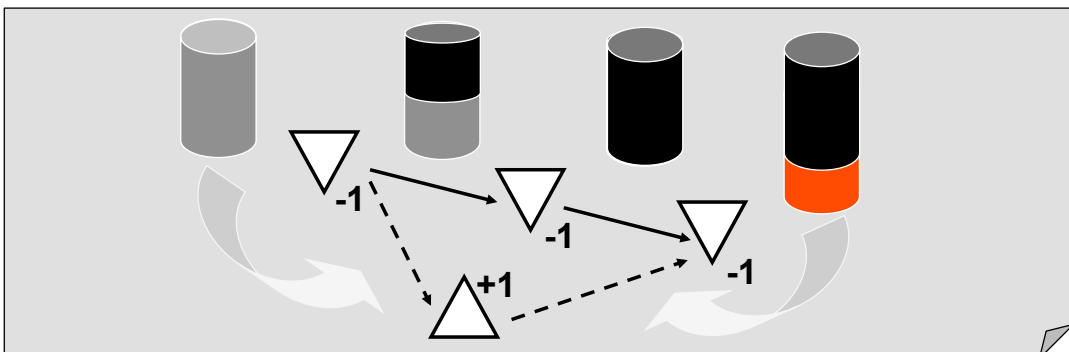
When time is relative (ordering of activities)

> then edge-finding and aggregated demand deduce nothing

We can still use information about ordering of activities and resource production/consumption!

**Example:**

> Reservoir: activities consume and supply items

Activity A „produces" **prod(A)** quantity:
  – positive number means **production**
  – negative number means **consumption**

- **Optimistic resource profile (orp)**
  – maximal possible level of the resource when A is processed
  – activities known to be before A are assumed together with the production activities that can be before A

$$orp(A) = InitLevel + prod(A) + \sum_{B<<A} prod(B) + \sum_{B??A \ \& \ prod(B)>0} prod(B)$$

- **Pessimistic resource profile (prp)**
  – minimal possible level of the resource when A is processed
  – activities known to be before A are assumed together with the consumption activities that can be before A

$$prp(A) = InitLevel + prod(A) + \sum_{B<<A} prod(B) + \sum_{B??A \ \& \ prod(B)<0} prod(B)$$

**\*B??A means that order of A and B is unknown yet**

$$orp(A) < MinLevel \Rightarrow fail$$
  – "despite the fact that all production is planned before A, the minimal required level in the resource is not reached"

$$orp(A) - prod(B) - \sum_{B<<C \ \& \ C??A \ \& \ prod(C)>0} prod(C) \ < MinLevel \Rightarrow B<<A,$$

for any B such that B??A and *prod*(B)>0
  – "if production in B is planned after A and the minimal required level in the resource is not reached then B must be before A"
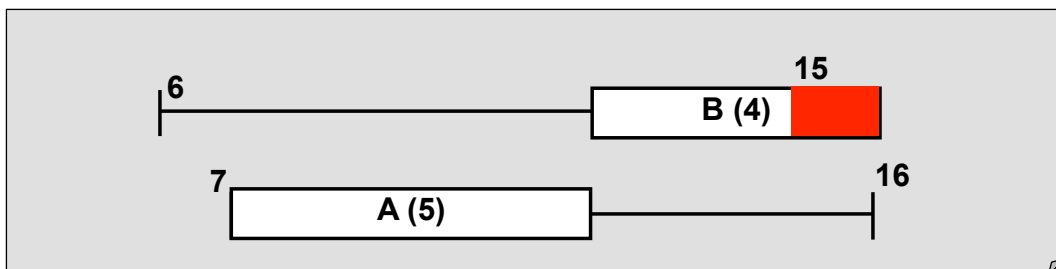
$prp(A) > MaxLevel \Rightarrow fail$

 – "despite the fact that all consumption is planned before A, the maximal required level (resource capacity) in the resource is exceeded"

$prp(A) - prod(B) - \sum_{B<<C \ \& \ C??A \ \& \ prod(C)<0} prod(C) > MaxLevel$
$\Rightarrow B<<A,$

for any B such that B??A and *prod*(B)<0

 – "if consumption in B is planned after A and the maximal required level in the resource is exceeded then B must be before A"

**What happens if A is processed before B?**



Restricted time windows can be used to infer inevitable precedence relations.

est(A)+p(A)+p(B)>lct(B) $\Rightarrow$ B«A

We can "energy" of activities processed before A to infer a more precise value of est(A)

**start(A) ≥ max{ est(Ω') + ⌈e(Ω')/cap⌉ | Ω'⊆{C : C«A} }**

## For unary resources:

**start(A) ≥ max{ est(Ω') + p(Ω') | Ω'⊆{C : C«A} }**
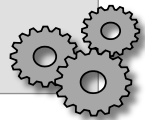
it is enough to explore sets $\Omega(X,A) = \{Y \mid Y « A \wedge est(X) \leq est(Y)\}$

**start(A) ≥ max{ est(Ω(X,A)) + p(Ω(X,A)) | X « A }**

```
dur ← 0
end ← est(A)
for each Y∈{ X | X « A } in the non-increasing order of est(Y) do
    dur ← dur + p(Y)
    end ← max(end, est(Y) + dur)
end for
est(A) ← end
```

The objective function in a CSP is usually encoded using an equality constraint and a new variable:

**v = obj(Xs)**

*Example:* **makespan = max{end(A$_i$)}**

– We can infer better bounds for variable v using the current domains of variables Xs:
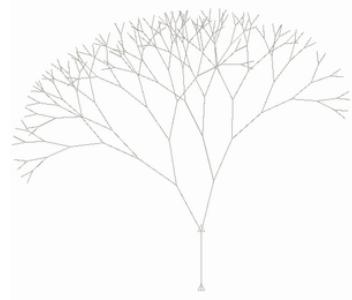
**makespan$_{min}$ = max{ect(A$_i$)}**

– We can infer better bounds for variables Xs using the current bounds for variable v:

**end(A$_i$) ≤ makespan$_{max}$**

– For more complex objective functions we can use problem relaxation to compute the bounds.

**Branching = resolving disjunctions**
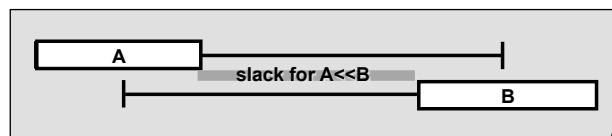
**Traditional scheduling approaches**:

- take **critical decisions first**
  - resolve bottlenecks …
  - defines the shape of the search tree
  - recall the **first-fail** principle
- prefer an **alternative leaving more flexibility**
  - defines order of branches to be explored
  - recall the **succeed-first** principle

How to describe criticality and flexibility formally?

**Slack** is a formal description of flexibility

- Slack for **a given order of two activities**
  „free time for shifting the activities"



slack(A<<B) = max(end(B))-min(start(A))-p({A,B})

- Slack for **two activities**
  slack({A,B}) = max{slack(A<<B),slack(B<<A)}

- Slack for **a group of activities**
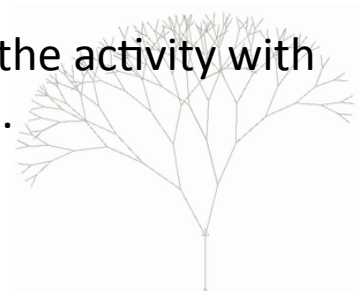  slack($\Omega$) = max(end($\Omega$)) - min(start($\Omega$)) - p($\Omega$)

**A<<B ∨ ¬A<<B**

- **Which activities should be ordered first?**
  - the most critical pair (first-fail)
  - the pair with **the minimal slack({A,B})**
- **What order should be selected?**
  - the most flexible order (succeed-first)
  - the order with **the maximal slack(A??B)**
- $O(n^2)$ choice points

**(A<<Ω ∨ ¬A<<Ω) or (Ω<<A ∨ ¬ Ω<<A)**

- **Should we look for first or last activity?**
  - select a **smaller set** among possible first or possible last activities (first-fail)
- **What activity should be selected?**
  - If first activity is being selected then the activity with the **smallest min(start(A))** is preferred.
  - If last activity is being selected then the activity with the **largest max(end(A))** is preferred.
- $O(n)$ choice points

**Resource slack** is defined as a slack of the set of activities processed by the resource.

## How to use a resource slack?

- choosing a resource on which **the activities will be ordered** first
  - resource with a minimal slack (**bottleneck**) preferred

- choosing a resource on which the **activity will be allocated**
  - resource with a maximal slack (**flexibility**) preferred