# Using constraint programming to resolve the multi-source/multi-site data movement paradigm on the Grid

Michal Zerola for the STAR collaboration
Jérôme Lauret, Roman Barták and Michal Šumbera

michal.zerola@ujf.cas.cz

NPI, AS CR

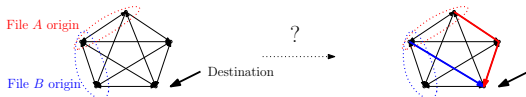ACAT 2008, November 4, 2008

# Outline of the problem

## Challenges

How to tranfer a set of files for a physicist to his site in a **shortest time** possible for fast analysis turn around?

How to achieve **controlled planning** of data-movement between sites on a grid (not necessarily following Tier architecture)?

Ultimately: How to **couple SE and CE** aspects for distributing data on the grid for further processing?

We concentrate on the first (most practical and of immediate need) question. On the grid:

- files are often replicated - available at multiple sites
- links diverse in bandwidth or QoS



We will present a Constraint Programming approach for solving this problem.

# Introduction to Constraint Programming

A Constraint Satisfaction Problem (CSP) consists of:

- a set of **variables** $X = \{x_1, ..., x_n\}$,
- for each variable $x_i$, a finite set $D_i$ of possible values (its **domain**),
- and a set of **constraints** restricting the values that the variables can simultaneously take.

A solution to a CSP is an assignment of a value from its domain to every variable, in such a way that every constraint is satisfied. We may want to find:

- any solution
- all solutions
- an optimal solution (defined by some objective function)

## Example - Sudoku

- logic-based number-placement puzzle
- the objective is to fill a $9 \times 9$ grid so that each column, each row, and each of the nine $3 \times 3$ boxes contains the digits from 1 to 9 **only one time each**
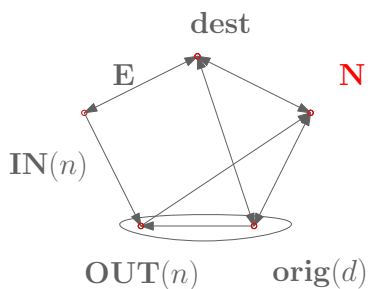


How to model Sudoku as a CSP?

- **variable** with a **domain** equal to $\{1, \ldots, 9\}$ for each single cell
- **constraint** of inequality between all pairs of variables in each row, column, or $3 \times 3$ cell
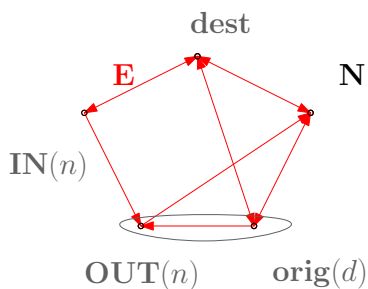
Notation:



**dest**

**E**

**N**

**IN**$(n)$

**OUT**$(n)$     **orig**$(d)$
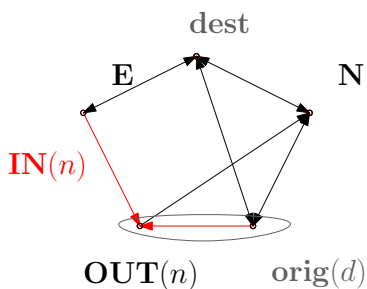
**D**- set of demands

# From Toy problem to the real life

Notation:



**D**- set of demands

Notation:



**D**- set of demands

Notation:



**D**- set of demands

Notation:



**D**- set of demands

Notation:



**dest**

**E**      **N**

**IN**$(n)$

**OUT**$(n)$      **orig**$(d)$

**D**- set of demands

Variables:

- Decision variable $\mathbf{X_{de}} = \begin{cases} 1, & d \text{ is routed over } e \\ 0, & d \text{ is not routed over } e \end{cases}$

- Integer variable $\mathbf{start_{de}}$ - start time of transfer $d$ over $e$
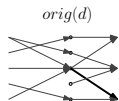
## Modeling constraints

Objective function - **minimizing makespan** (latest finish time):

$$\min_{X_{de}, start_{de}} \max_{e \in \mathbf{E}} \underbrace{\left( start_{de} + \frac{size(d)}{speed(e)} \right)}_{end_{de}} \cdot X_{de}$$
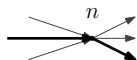
Extraction of constraints:

- for each file data transfer must start from exactly one of its origins



$$\forall d \in \mathbf{D} : \sum_{e \in \cup \mathbf{OUT}(n|n \in \mathbf{orig}(d))} X_{de} = 1, \sum_{e \in \cup \mathbf{IN}(n|n \in \mathbf{orig}(d))} X_{de} = 0$$

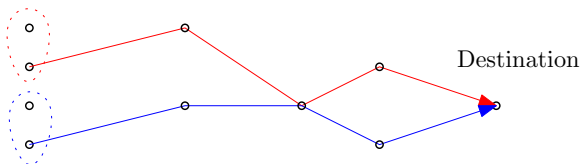- if the file is transferred trough a given site, it must enter and leave site exactly once



$$\forall d \in \mathbf{D}, \forall n \notin \{orig(d) \cup dest(d)\} :$$
$$\sum_{e \in \mathbf{OUT}(n)} X_{de} \leq 1, \sum_{e \in \mathbf{IN}(n)} X_{de} \leq 1, \sum_{e \in \mathbf{OUT}(n)} X_{de} = \sum_{e \in \mathbf{IN}(n)} X_{de}$$

- $\cdots$ and many others

# Solution finding principle

The solving is composed of tree search-like **iterations over two stages**:

1. choose a path for each demand from its origins to the destination (assigning $X$ variables)



Destination

2. for each link schedule all transfers, as on disjunctive resource (assigning *start* variables)

## Complexity

The second stage corresponds to job-shop scheduling that belongs to the $NP - hard$ class of problems.

# Pruning the search space

We cannot escape exponential time growth but we can do the best to shift it, so our input instances will be solved fast.
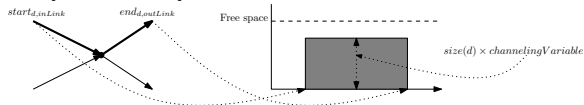
How to achieve it?

- **symmetry breaking** - don't explore configurations that are symmetrical (e.g. require ordering of the files from the same origin)

- **branch cutting** - don't explore transfer paths that cannot lead to the better schedule than we already have (too many file transfers per some link can make a lower bound estimate of the makespan)

- **heuristic** can lead to a fast (not necessarily optimum) solution, but "good enough"

... and other techniques, open for research
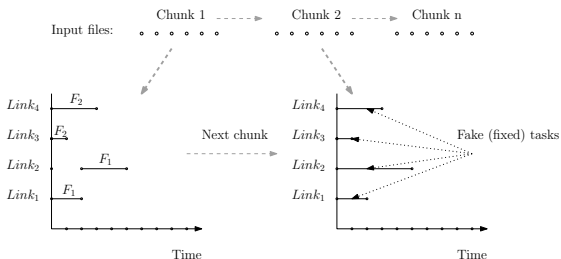
# Additional real life constraints

- sites have a **limited storage space**
  - easily modeled by *cumulative* constraint for each resource at a site



- several {in,out}coming links can have a **shared** bandwidth
  - easily modeled by *cumulative* constraint for each resource at a shared link
- **fair share** in a multiuser environment
  - scheduling by chunks, fair share corresponds to the selecting approach from waited demands
- network characteristic is not static but **fluctuates** in time
  - scheduling by chunks, sensors

# Scheduling by chunks

We can split scheduling of the whole file set into parts (chunks) and create
an optimal schedule for each part separately, while propagating previous
results.



Benefits:

- speed
- self adaptation to the network
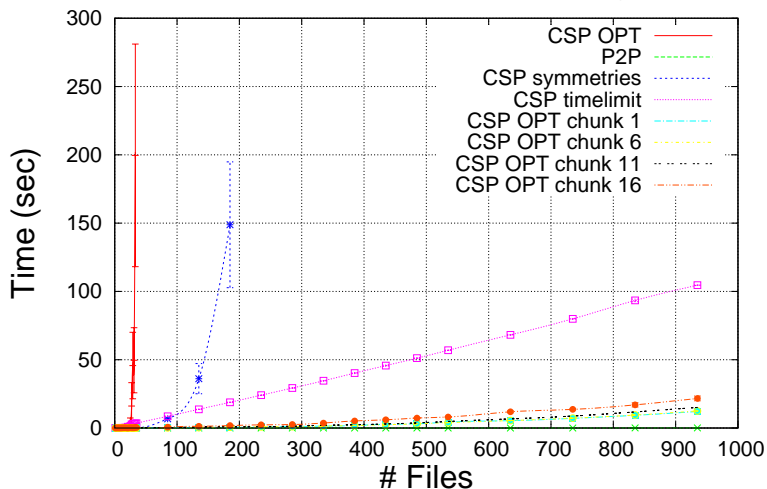- fair-share (responsiveness)

## Implementation

We choosed Java library for implementation of the scheduler. The main heart of the solver is inside an open-source **Choco** library designed for CSP and CP.
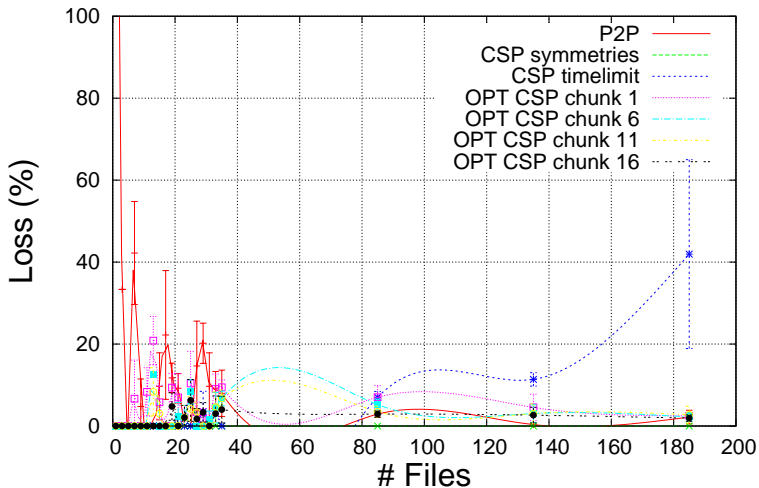
http://sourceforge.net/projects/choco/

- we implemented a simulator scheduling file transfers generated by feeders
- feeder generates file distribution with file origins:
    - distinct (file is available only on 1 site)
    - weighted (realistic probability function)
    - shared (file is available on each site)
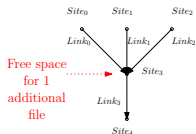- performance is compared to the Peer-2-Peer system

Time to produce a schedule (weighted)

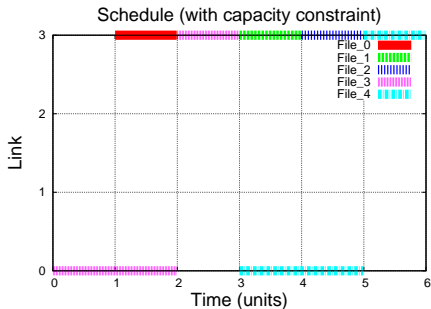# Results (2)
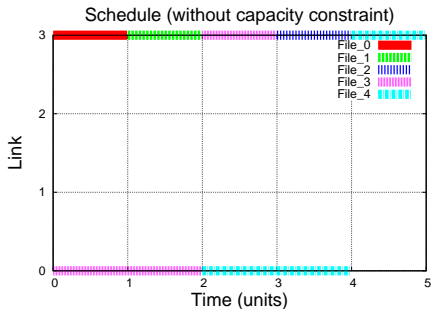


Makespan loss on optimum (weighted)

# Effect of the capacity constraint



### Gantt charts:

## Conclusions and future work

Constraint programming seem most suitable as:

- it allows to express in human terms real life problems
- it has a declarative character that allows easier further extensions

Further work would need to include:

- Enhancements for the $2^{nd}$ stage scheduling phase implementing algorithms from the job-shop problem field
- Testing several heuristics
- Move from simulation to the real life environment

Our first tests show promising results for single destination (site) data-movements.

- Out-performs p2p
- Performs well with planning-by-pieces (chunks)
- Hence, scales to large number of files