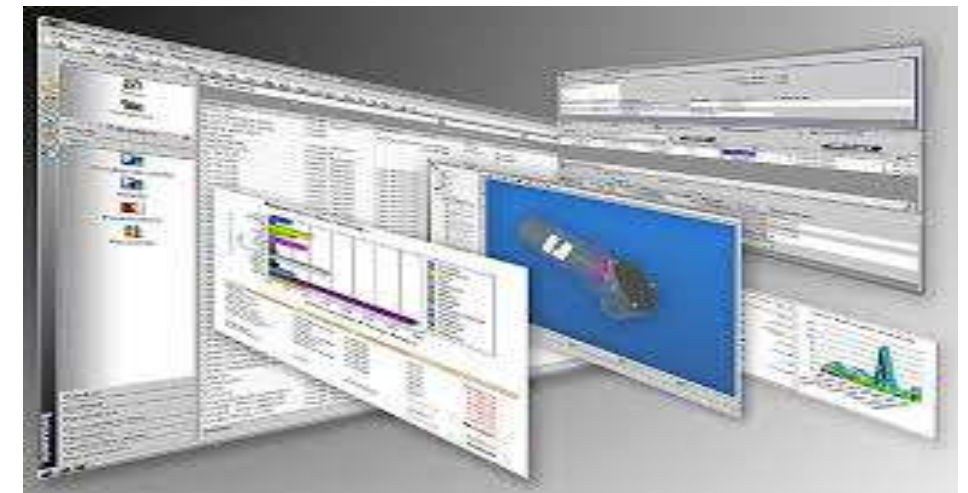*Hisham Ihshaish*

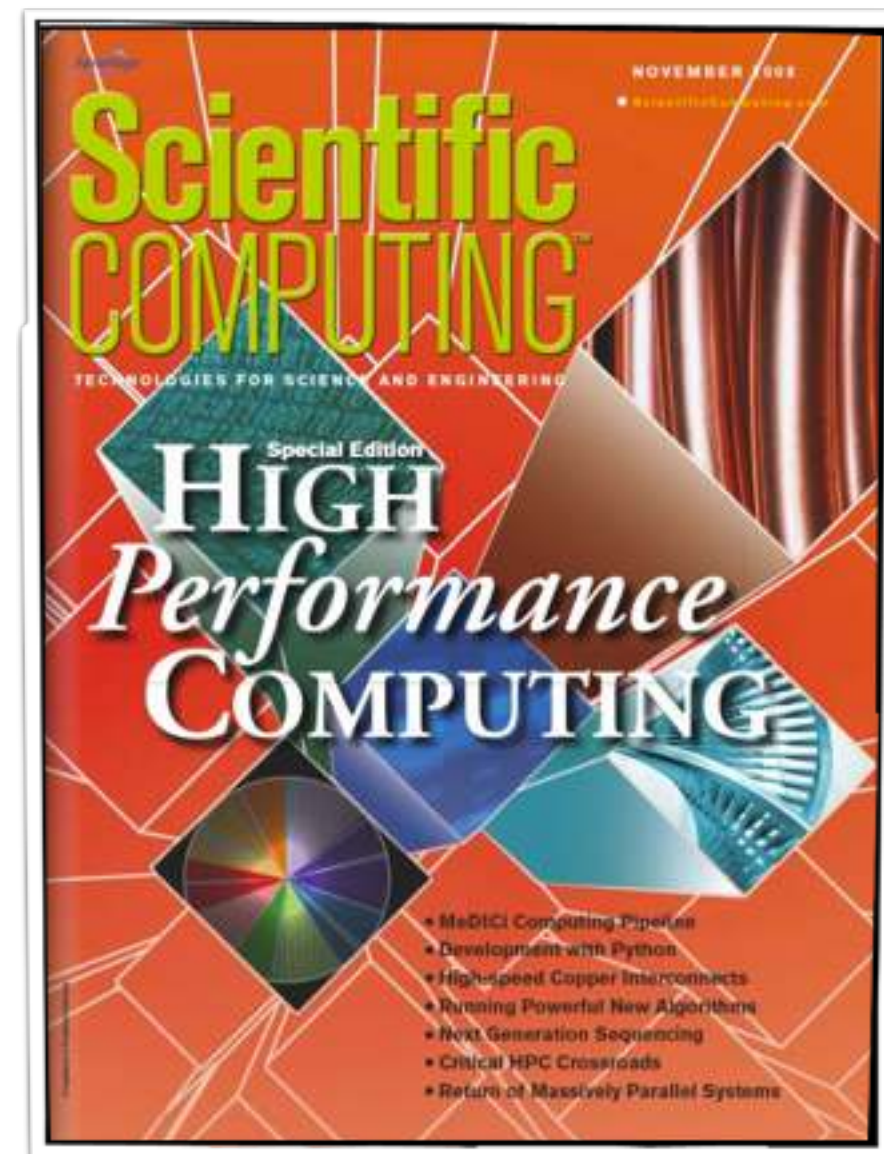# Parallel Software Tools for the Construction and Analysis of Complex Networks



LINC Workshop 4 – Montevideo
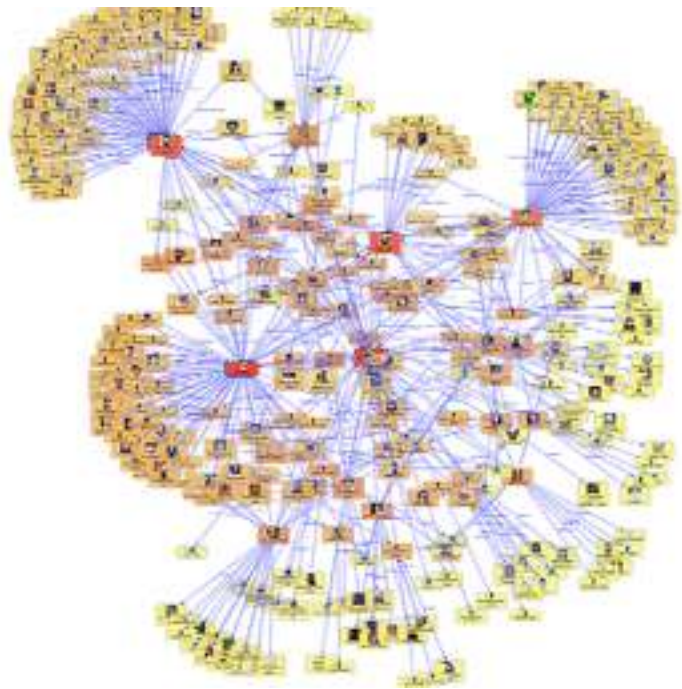24-26/March/2014

# Outline

- ❖ Introduction
- ❖ Computational challenges
- ❖ Parallel Software tools
    - ❖ Network construction
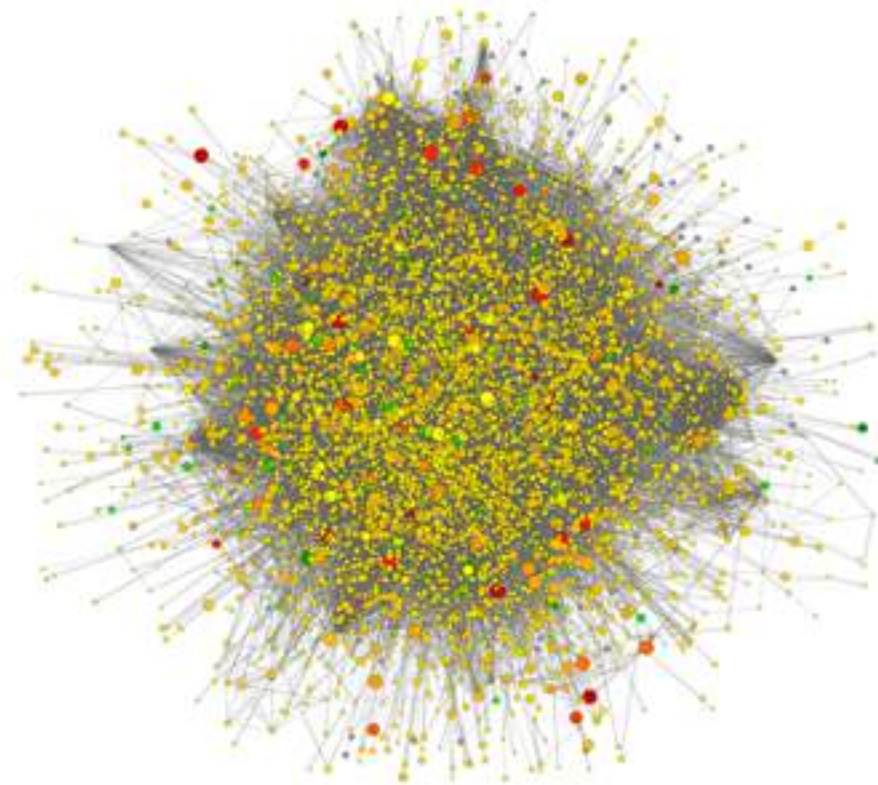    - ❖ Network analysis
- ❖ Software architecture
- ❖ Usage

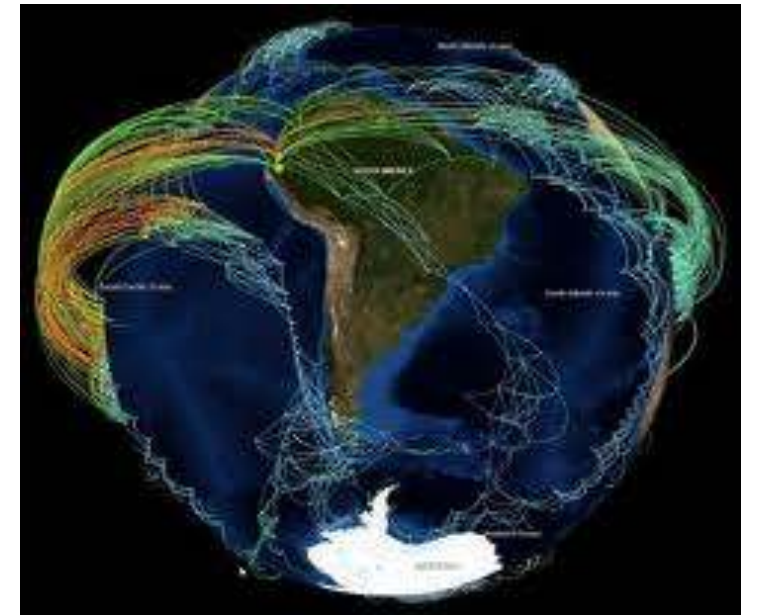# Introduction – large-scale data analysis

## Social Networks



Application: e.g., identifying communities, information spread modeling

## Bioinformatics



Application: e.g., identifying drug target proteins

## Climate



Application: e.g., identify patterns? analyze spatio-temporal interaction of climate variables

➡ **Sources of data:** simulations, experimental devices, the Internet, sensor networks

➡ **Challenges:** data size, heterogeneity, uncertainty, data quality, computational time

# Introduction – from domain-specific to computation

## application

Social Network Analysis

WWW

Computational Biology

Scientific Computing

Climate Research

## problems

community detection, central entities

marketing
social search

metabolic pathways, gene regulation

graph partitioning, coloring, matching

community detection, teleconnections

**data size** →

**complexity** →

## graph algorithms ✚

traversals, shortest paths

centrality measures

connectivity

community detection

## architecture

• Single processing unit

• Parallel machines

- GPUs
- x86 multicore servers
- Massively multithreaded clusters, ....
- Multicore clusters,
- Distributed memory clusters
- Clouds

# Introduction – from domain-specific to computation

**Input data**

⬇

## Network

⬇

*find ..*

- paths
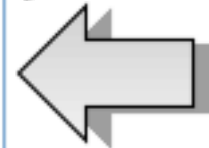- clusters
- partitions
- matchings
- patterns
- orderings

➡

## Graph kernel

- traversal
- shortest path algorithms
- flow algorithms
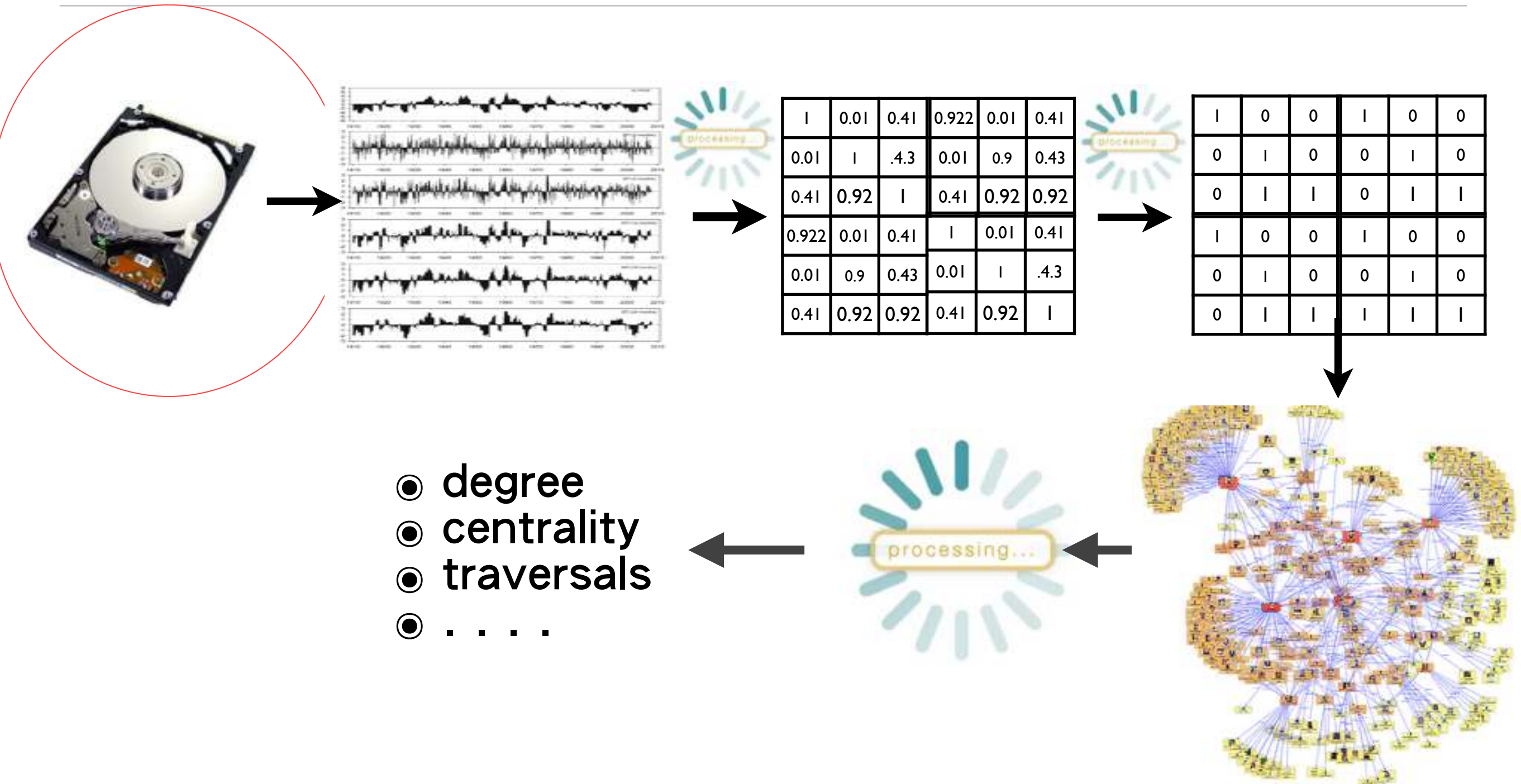- spanning tree algorithms
- topological sort
…..

⬅

**Which algorithm? factors.....**

- graph sparsity
- static/dynamic nature
- weighted/unweighted, weight distribution
- vertex degree distribution
- directed/undirected
- simple/multi/hyper graph
- problem size
- granularity of computation at nodes/edges
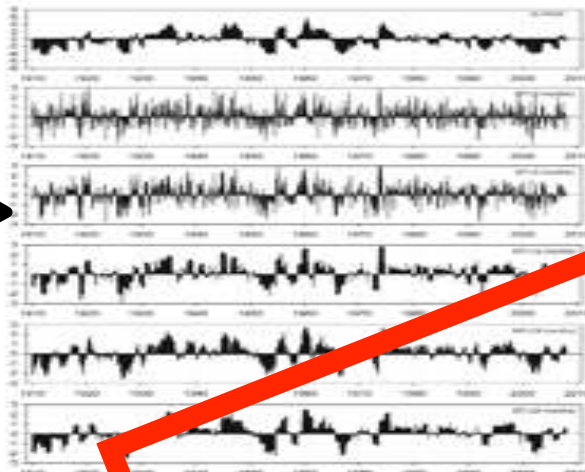- domain-specific characteristics

⬇ ⬆

## Computing architecture

# Introduction – computational challenges



- ◉ degree
- ◉ centrality
- ◉ traversals
- ◉ . . . .

# Hard "drive"

# Introduction – computational challenges



SLOW

BIG

Slow

- ◉ degree
- ◉ centrality
- ◉ traversals
- ◉ . . . .

# Introduction – computational challenges
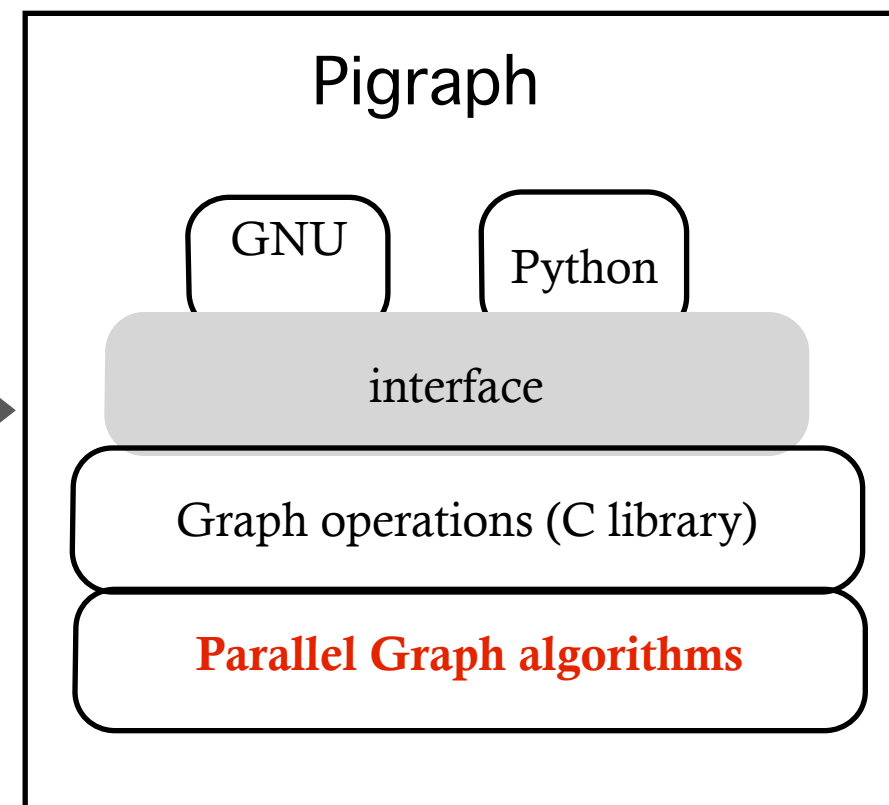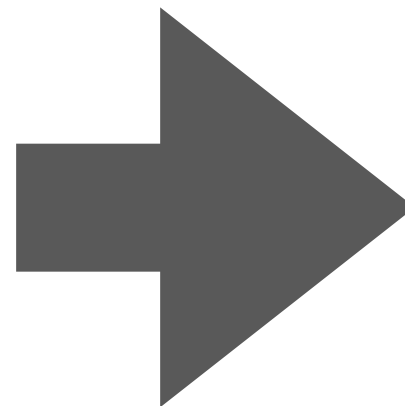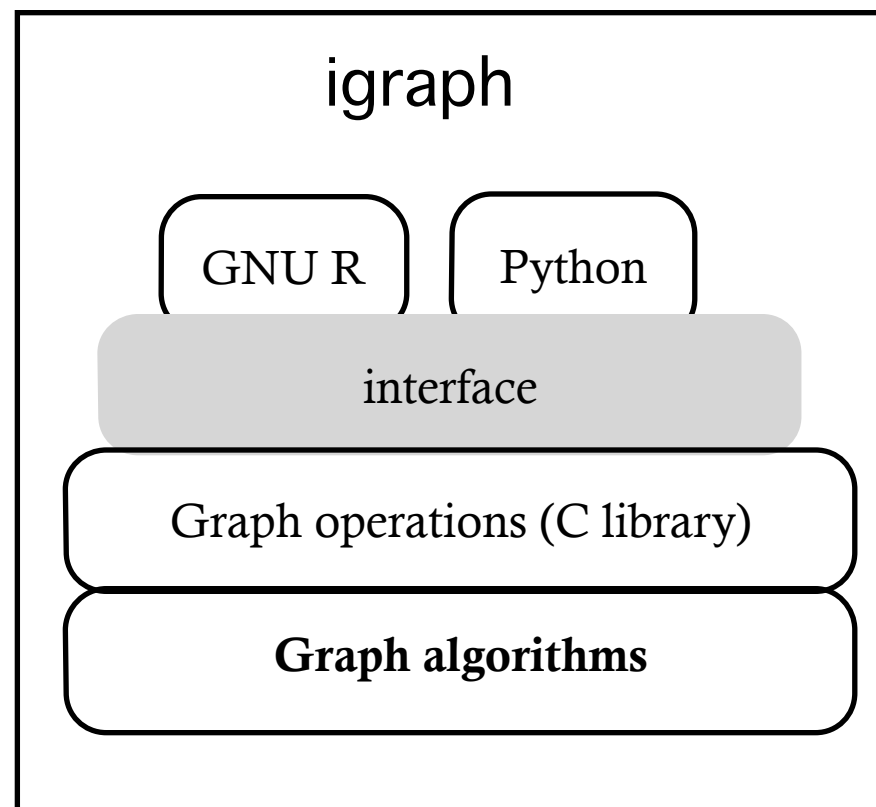


SLOW

BIG

SLOW

Slow

- ◉ parallel reading?
- ◉ data reduction/compression?
- ◉ parallel tools to analyze complex networks?
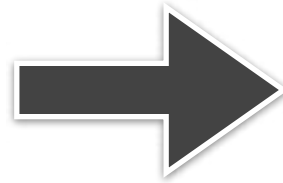
Parallel tools to analyze complex networks

Pigraph

# Pigraph – parallel library for graph analysis



➡ OpenMP

➡ Shared memory platforms

# Pigraph – example: "shortest path": Floyd-Warshall algorithm

$$D(0) = \begin{bmatrix} 0 & 3 & 8 & \infty & -4 \\ \infty & 0 & \infty & 1 & 7 \\ \infty & 4 & 0 & \infty & \infty \\ 2 & \infty & -5 & 0 & \infty \end{bmatrix}$$
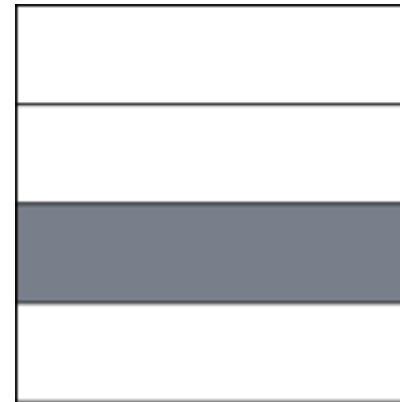
$$D_{(k-1)} = \begin{bmatrix} 0 & 3 & 8 & \infty & -4 \\ \infty & 0 & \infty & 1 & 7 \\ \infty & 4 & 0 & \infty & \infty \\ 2 & 5 & -5 & 0 & -2 \\ \infty & \infty & \infty & 6 & 0 \end{bmatrix}$$

```
For  k=1 to n {            Parallel
  For  i=1 to n {
     For  j=1 to n
        D[i,j] = min(D[i,j],D[i,k]+D[k,j])
  }
}
```
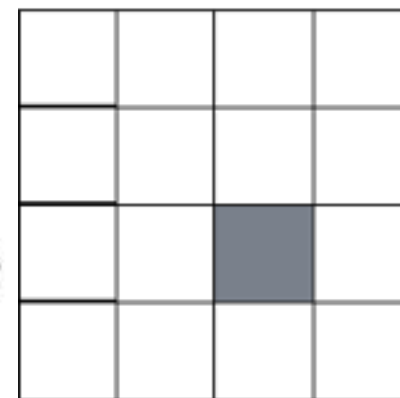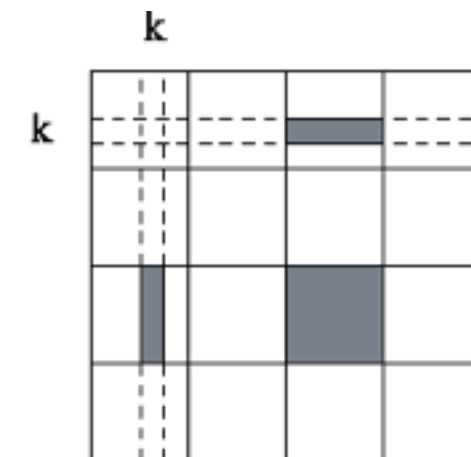


(a)          (b)

```
For  k=1 to n {            Parallel
  For  i=1 to n {          Parallel
     For  j=1 to n
        D[i,j] = min(D[i,j],D[i,k]+D[k,j])
  }
}
```
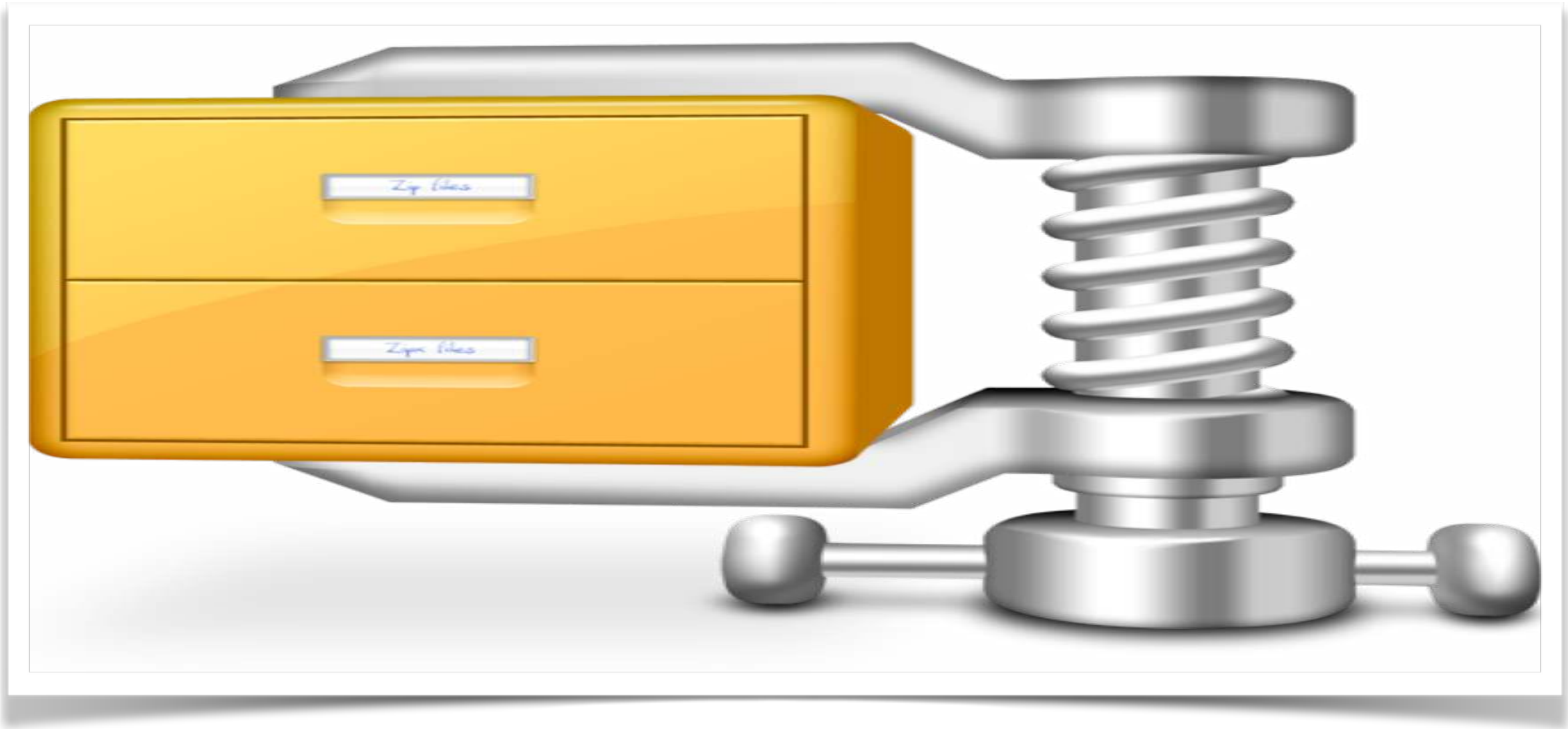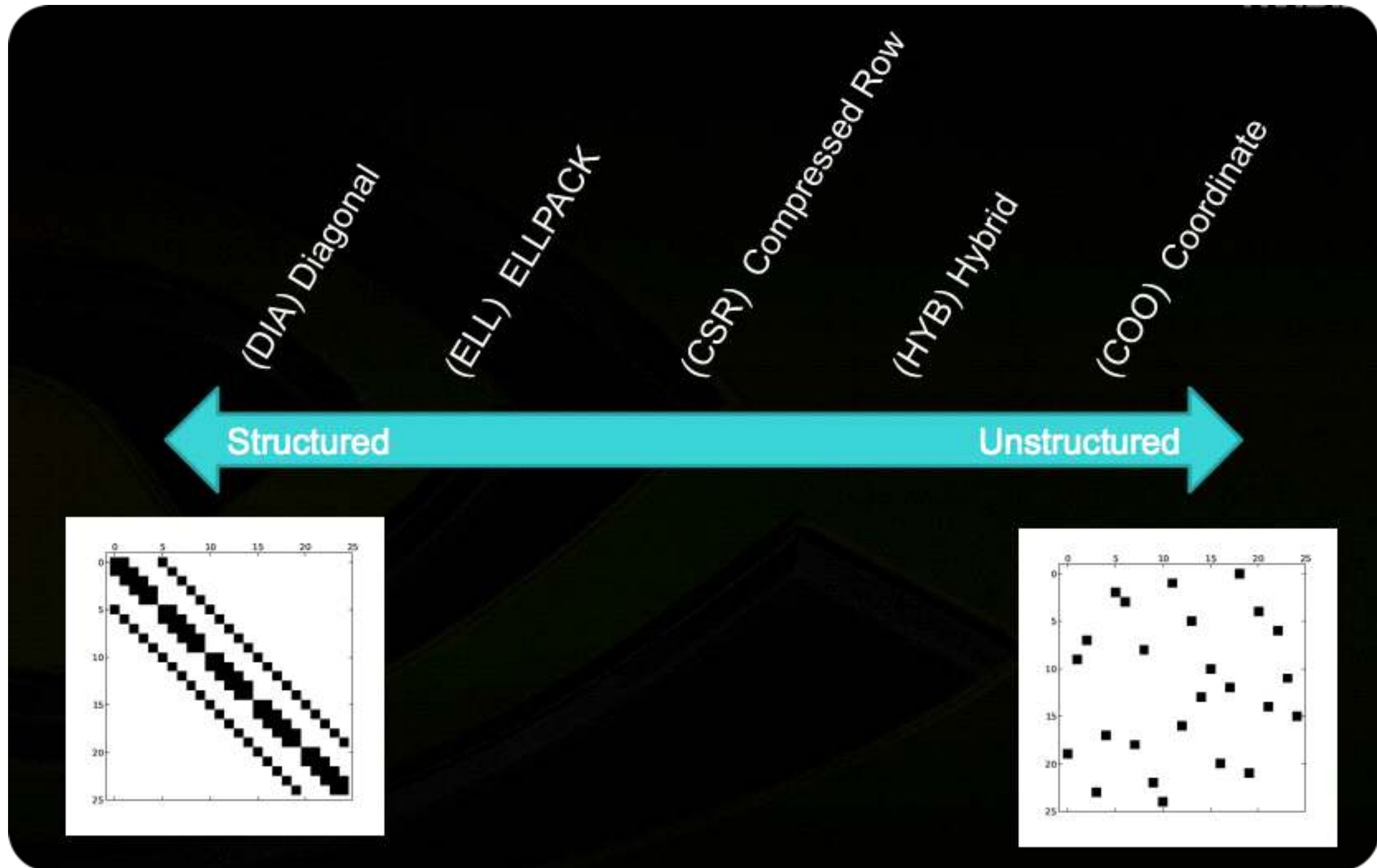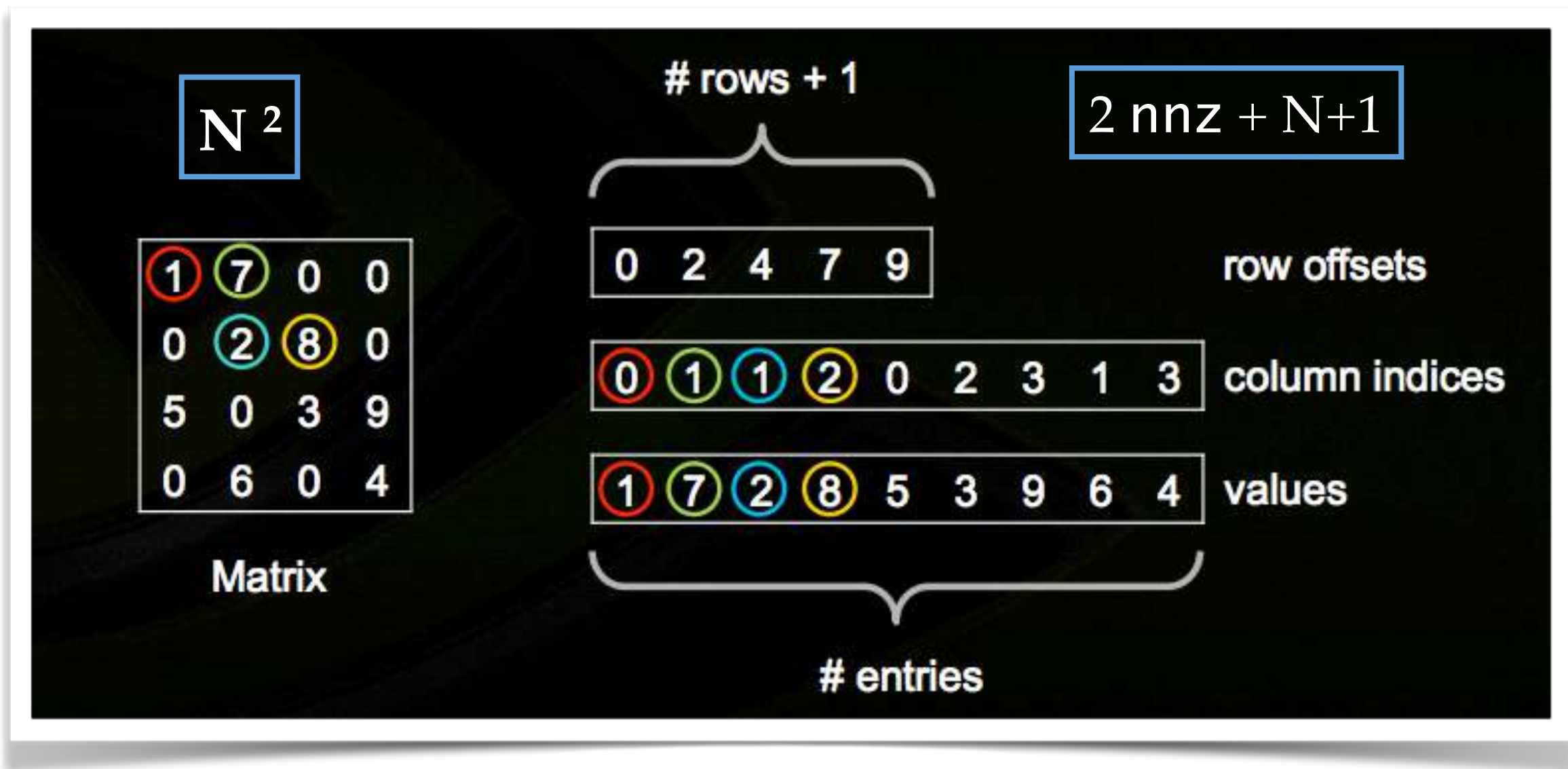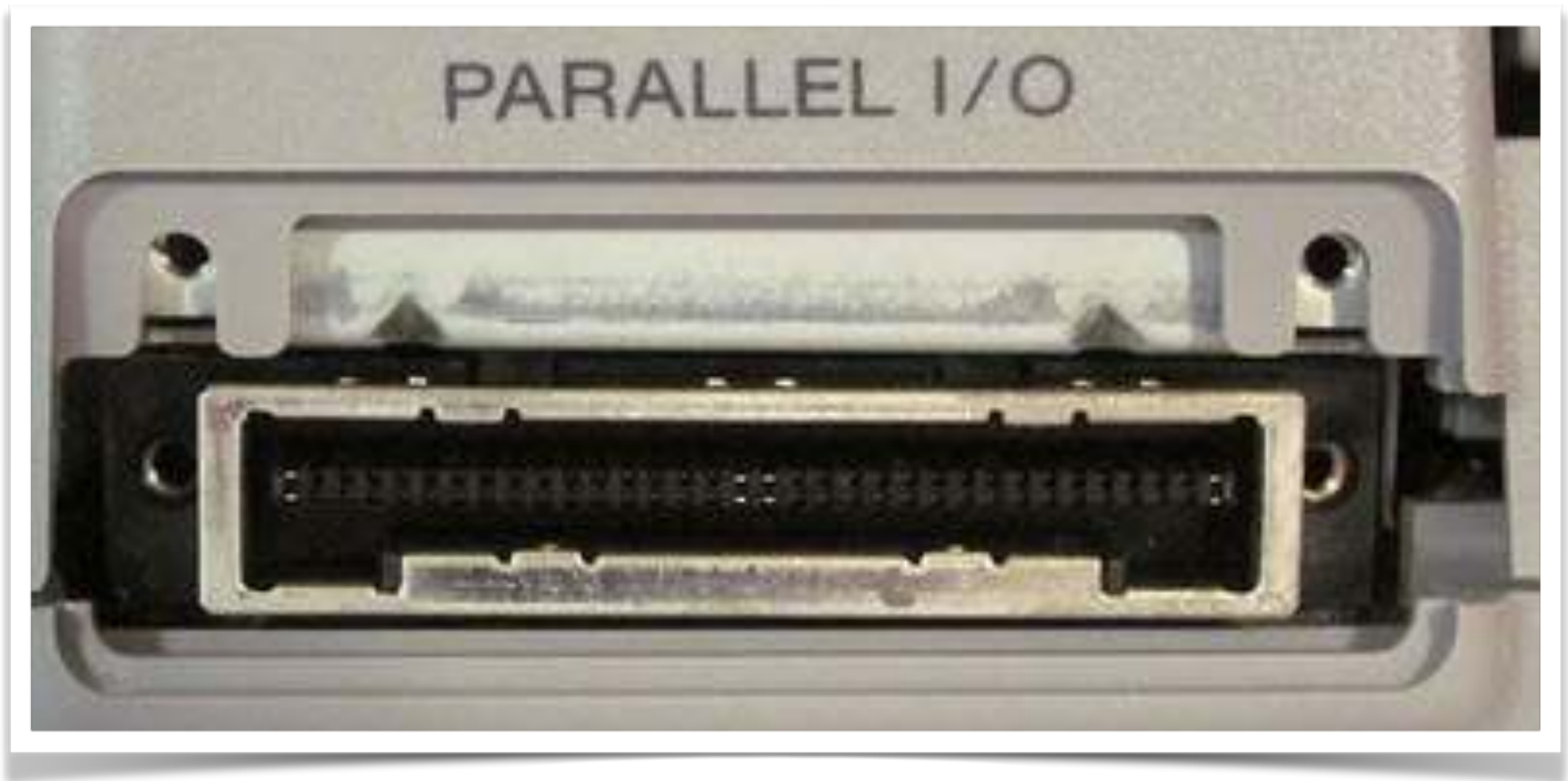


(a)          (b)

Data reduction / compression          CSR

# Compressed sparse matrices

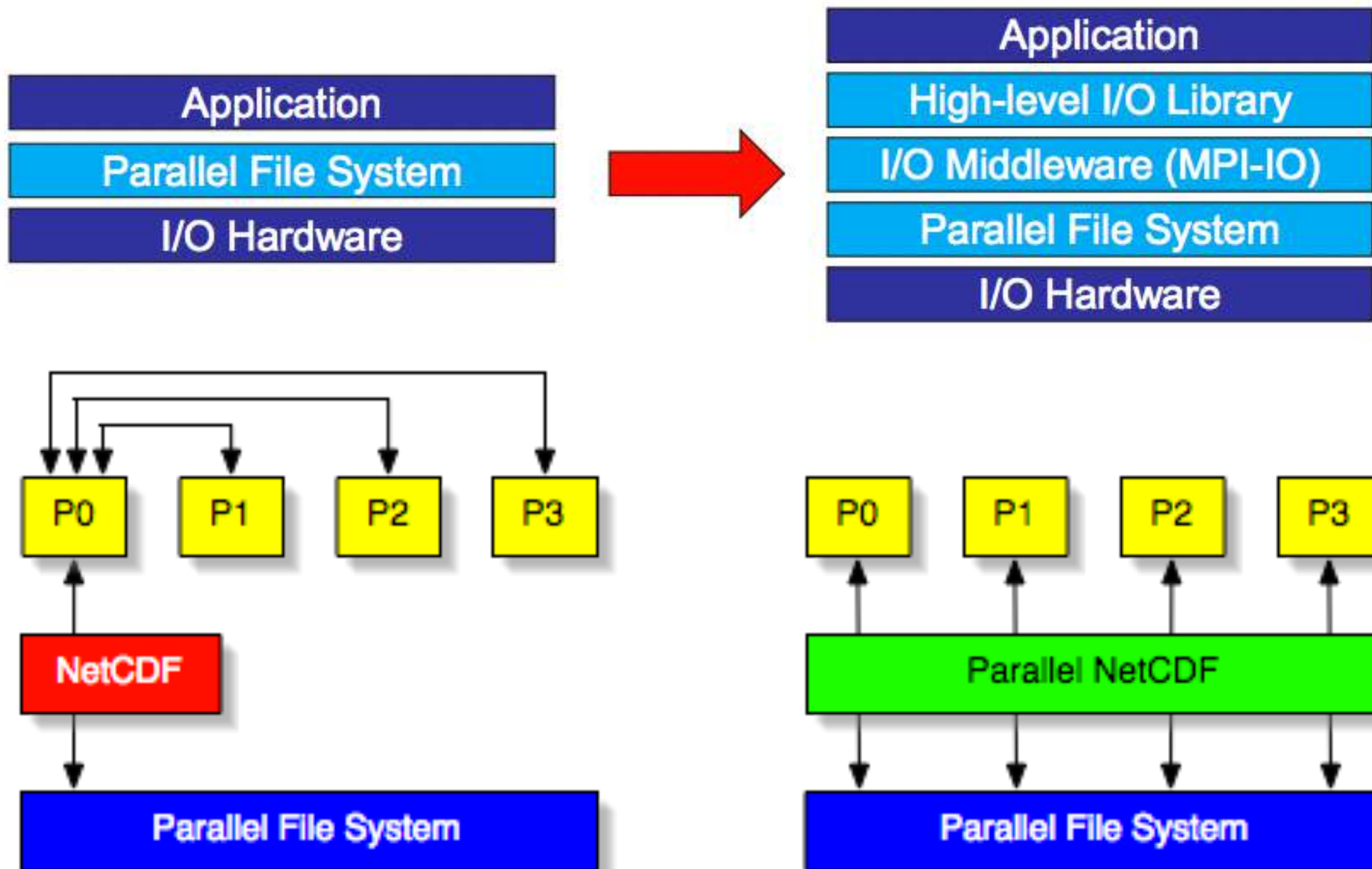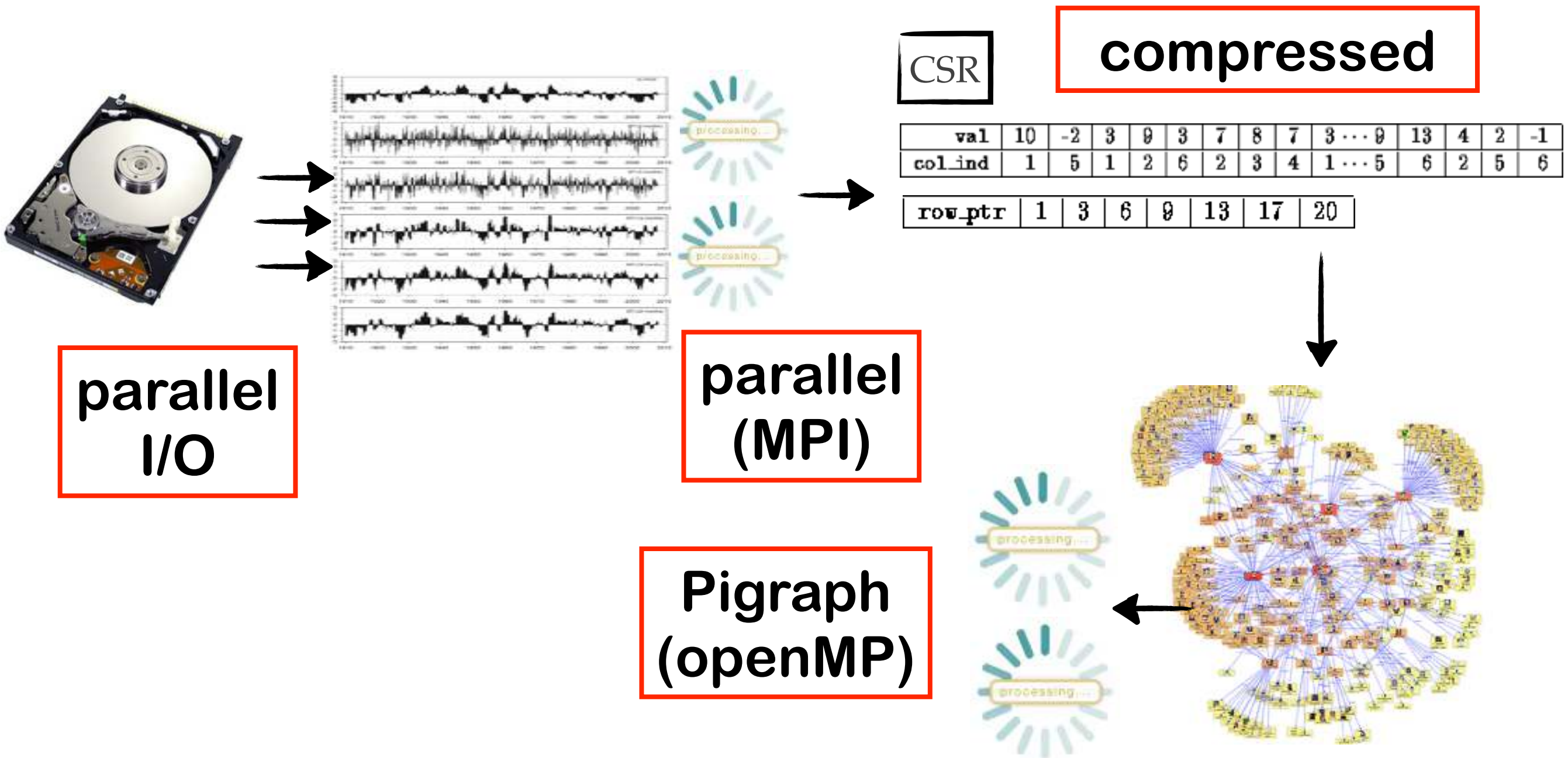# Compressed sparse matrices

**Compressed Sparse Row (CSR)**

Parallel reading

NetCDF-4
(HDF5, MPI/IO)

# Parallel I/O

# Parallel software tools for the construction and analysis of complex networks



**compressed**

**CSR**

| val | 10 | -2 | 3 | 9 | 3 | 7 | 8 | 7 | 3 ⋯ 9 | 13 | 4 | 2 | -1 |
|-----|----|----|---|---|---|---|---|---|-------|----|---|---|----|
| col_ind | 1 | 5 | 1 | 2 | 6 | 2 | 3 | 4 | 1 ⋯ 5 | 6 | 2 | 5 | 6 |

| row_ptr | 1 | 3 | 6 | 9 | 13 | 17 | 20 |
|---------|---|---|---|---|----|----|----|

**parallel I/O**

**parallel (MPI)**

**Pigraph (openMP)**

# Parallel software tools for the construction and analysis of complex networks

# Thanks!