

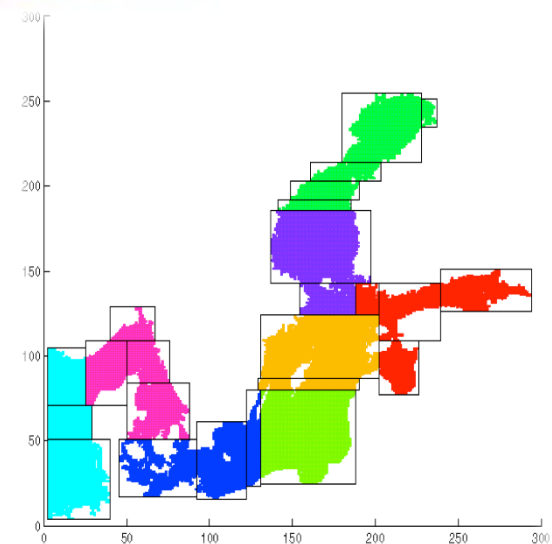
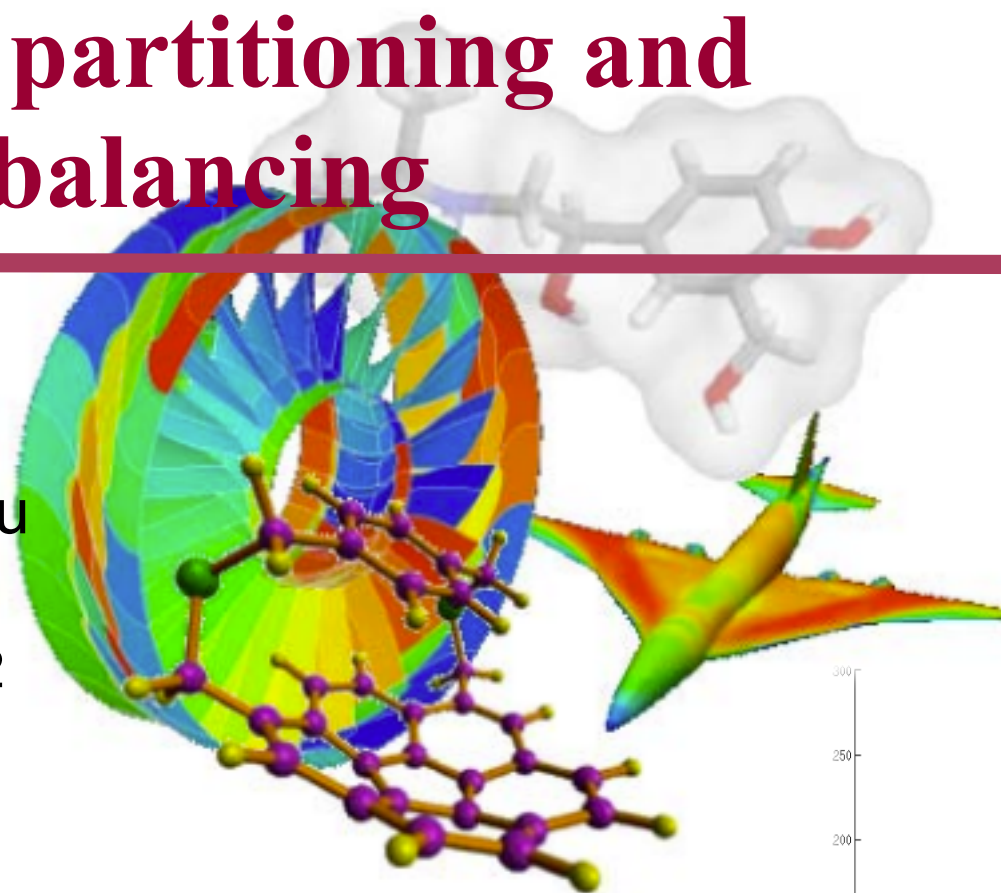


UPPSALA  
UNIVERSITET

# Data partitioning and load balancing

Jing Liu

2016-02





UPPSALA  
UNIVERSITET





## Case study:

Consider 7 independent tasks with estimated work loads: 5, 2, 3, 4, 5, 2, 10

Assign the tasks to 2 processors and compute the work load for each processor.

(How would you do with 10000 tasks and 100 processors?)

# Linear partitioning (array of N tasks)

## 1. *Static partitioning*

set  $n = \text{floor}(N/p)$ ,  $m = \text{mod}(N/p)$

Assign the  $m$  first proc  $n+1$  tasks  
and the other  $n$  tasks

Ex.  $N=10$ ,  $p=4 \Rightarrow n=2$ ,  $m=2$

tasks: 3,3,2,2



- + very easy
- without considering workload of each task

Example: operations on dense matrix

# Linear partitioning (array of N tasks)

## 2. *Cyclic partitioning*

Assign tasks cyclicly (Round Robin)

$\text{proc}(i) \leq \text{task}(i+n*p), n=1,2,\dots$

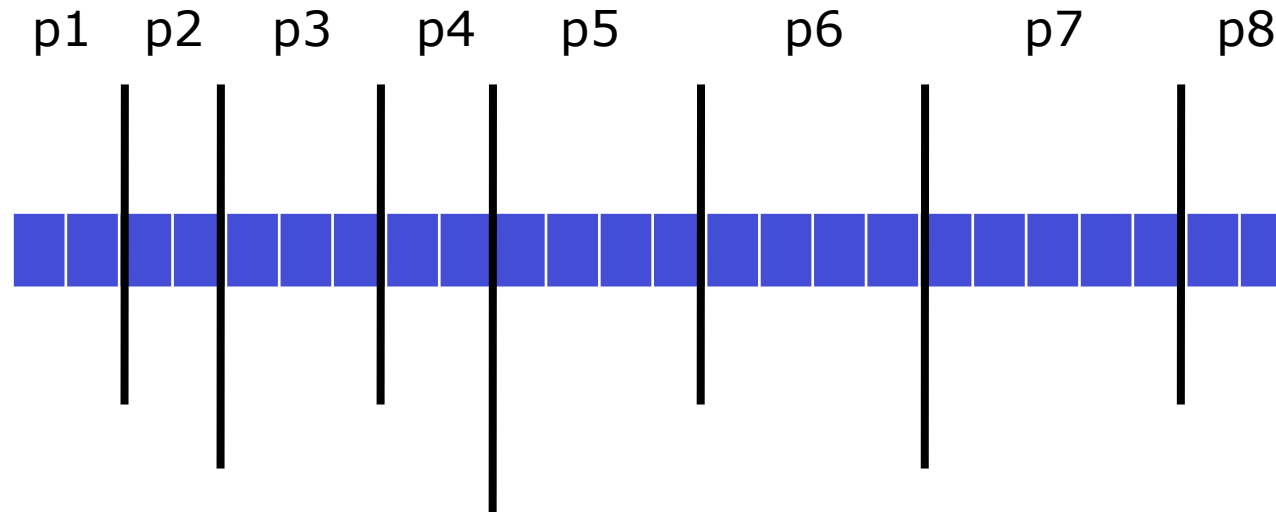


- + easy
- without considering workload of each task, but since Round Robin is used, we may be less imbalanced.

ex: LU-factorization, Gram-Schmidt, Gaussian elimination

### 3. *Recursive bisection*

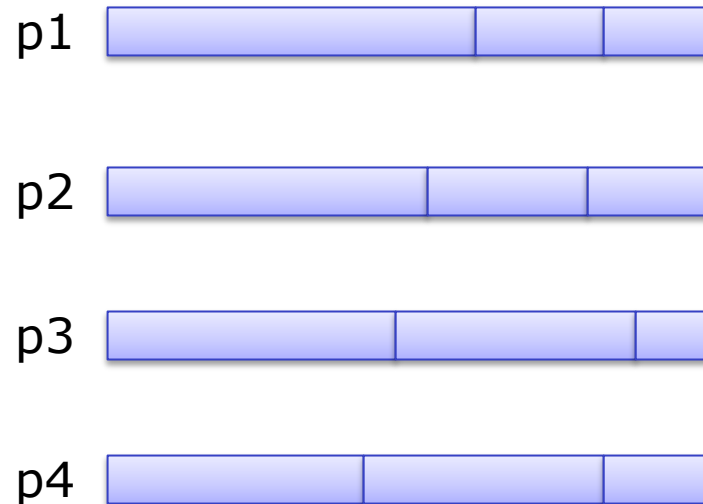
Divide the array into two halves with equal work loads. Proceed recursively with each half until we have  $p$  parts.



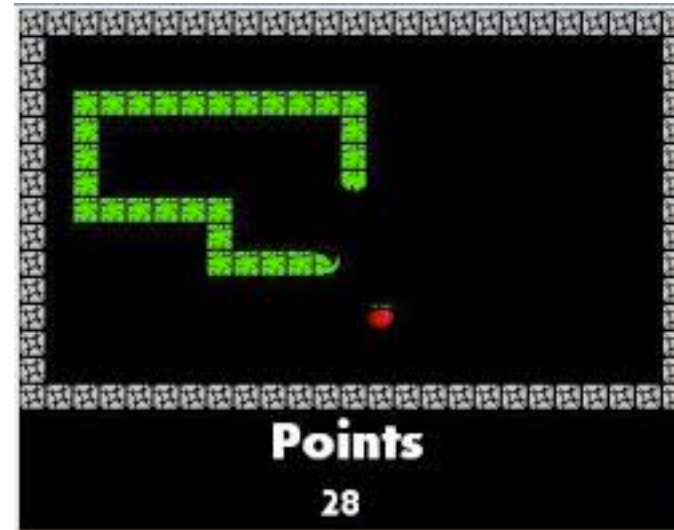
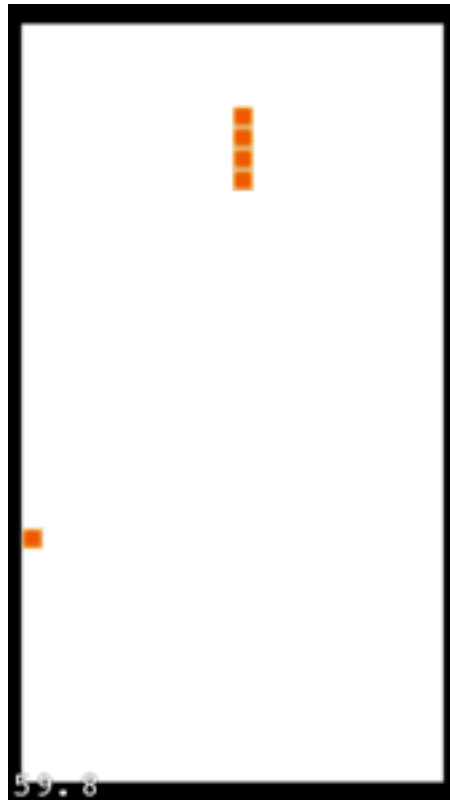
Can easily be extended to arbitrary number of processors by assigning the work load proportional to the processors in each cut.

## 4. *Bin-packing*

Assign tasks in size order, **largest** first, to the processor with least work in each step.



## 5. Greedy (variant of Bin-packing)







## 5. Greedy (variant of Bin-packing)

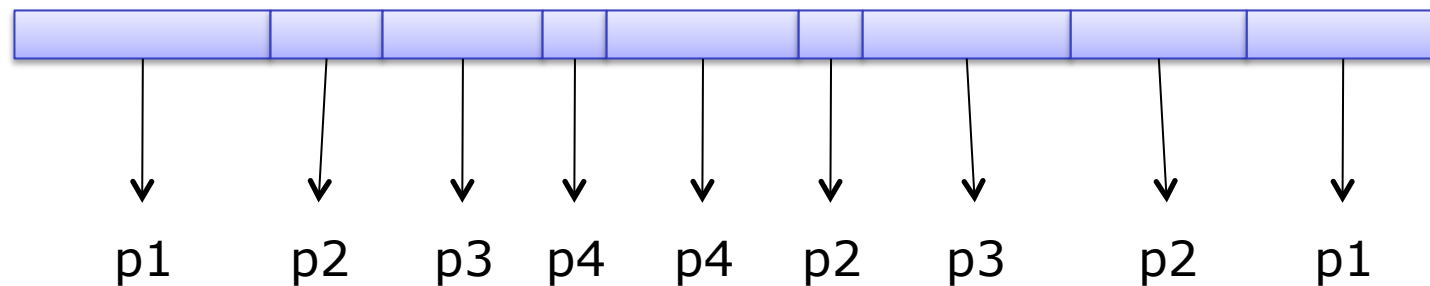
1. Set an upper limit  $C$  for the work load in each bin.
2. For each bin grab as many tasks as possible (load  $\leq C$ ).
3. If we could not assign all tasks, increase  $C$  otherwise decrease  $C$ .

Repeat 1-3 searching for the smallest possible  $C$  for which it is possible to assign all tasks.

(Lower complexity than for Bin-packing)

## 6. *Dynamic*

Set up a task-queue and assign tasks from the queue to the processors as soon as they are ready processing a previous task.



(See example 7.4 and 7.6 for Pthreads implementation, OpenMP supports `schedule(dynamic)` and task queues.)

**Case study:** Tasks: 5, 2, 3, 4, 5, 2, 10

Static:  $5+2+3+4=14$ ,  
 $5+2+10=17$

Cyclic:  $5+3+5+10=23$   
 $2+4+2=8$

Bisect:  $5+2+3+4=14$   
 $5+2+10=17$

Bin-pack:  $10+4+2=16$   
 $5+5+3+2=15$

Greedy:  $10+4+2=16$   
 $5+5+3+2=15$

Dynamic:  $5+4+2=11$   
 $2+3+5+10=20$

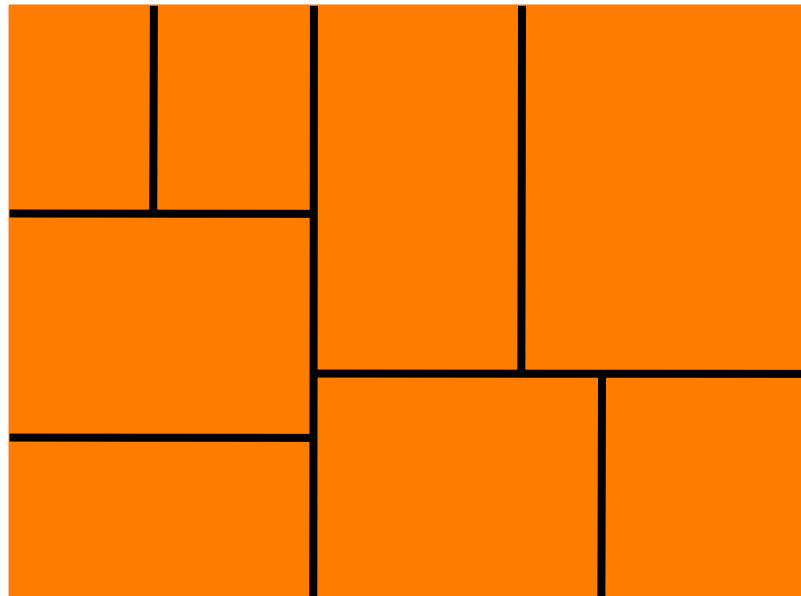
# Data partitioning in 1D

- Non-Weighted
  - all tasks are equal
  - only considers the number of processors and tasks
  - easy and fast
- Weighted
  - each task has a weight
  - more complex, takes more time
  - may give a better performance

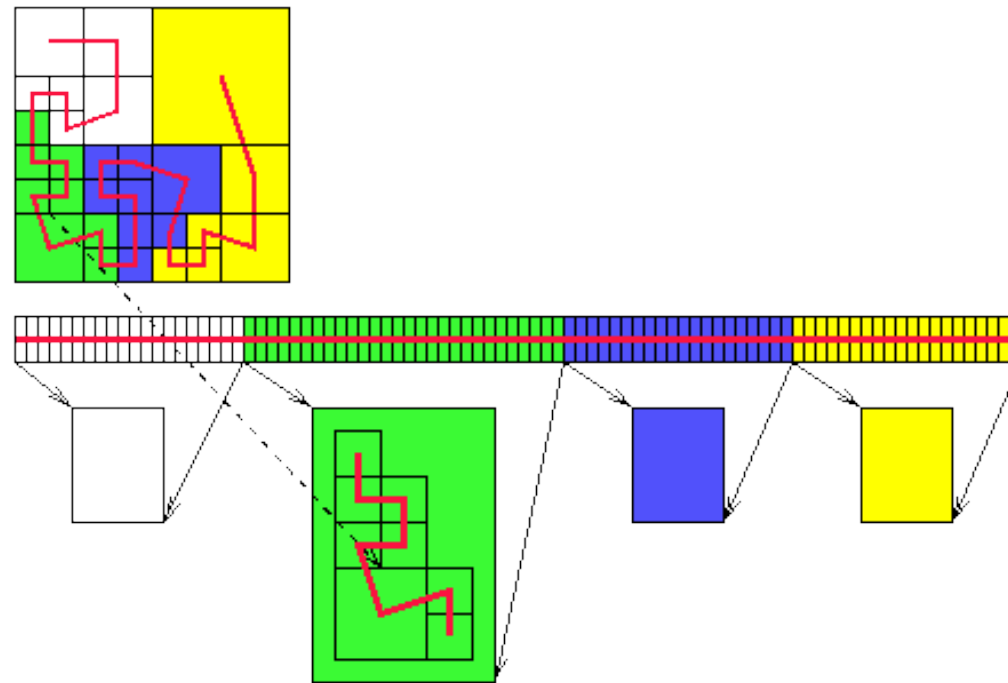
# Data partitioning in n-D

**Note:** The algorithms can easily be extended to multi-dimensional arrays (for partitioning grids or matrices).

Ex Recursive bisection:



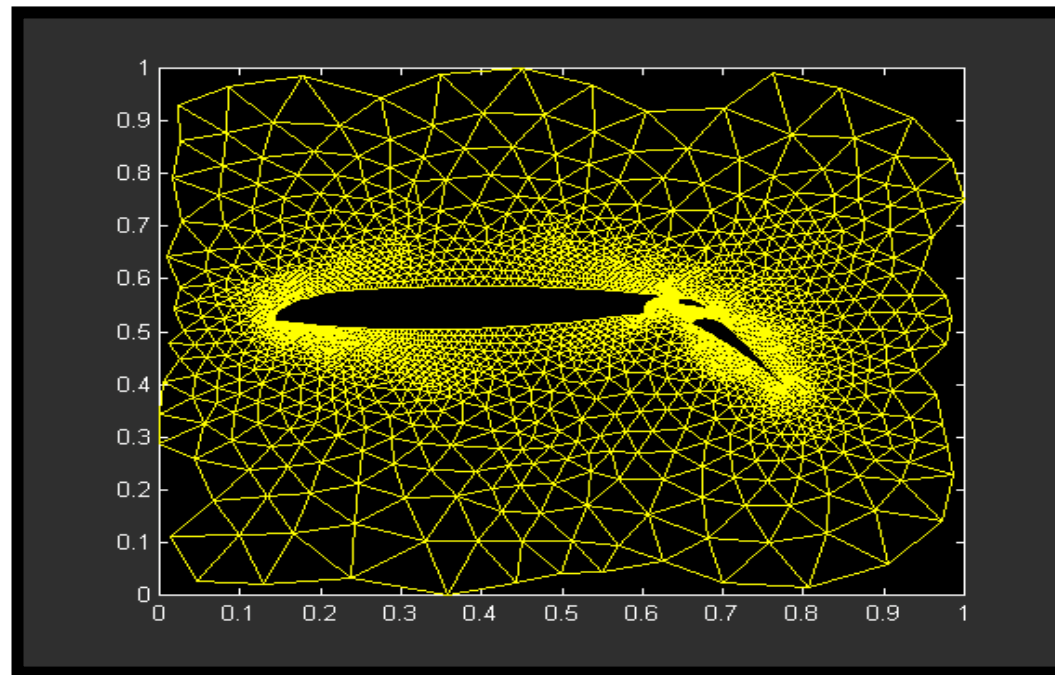
Or by mapping to a one dimensional array with Space Filling Curves (SFC), e.g., Hilbert or Morton curve.



SFCs has a locality preserving property giving geometrically collected partitions.

# Graph partitioning methods

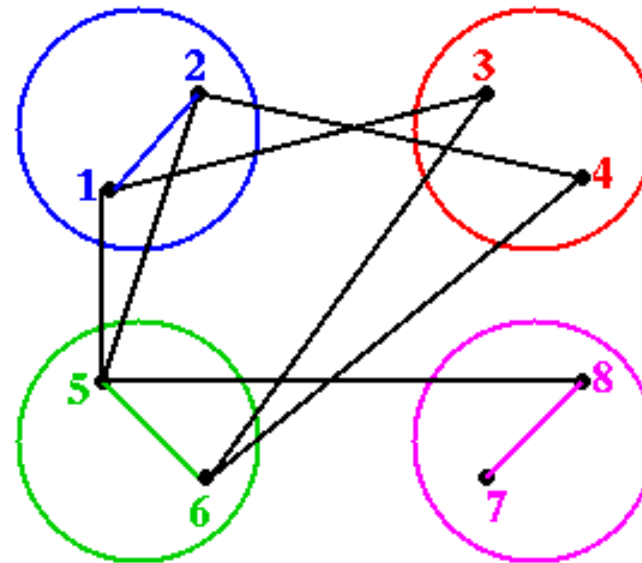
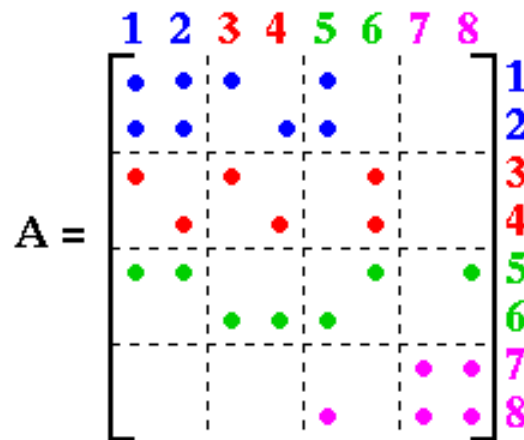
Assume that we have data with dependencies to nearest neighbors. These dependencies can then be represented with a graph, e.g., consider a FEM-mesh or sparse matrix mult.



Nodes - data, Edges - data dependencies

We want to partition the data equally (load balance) with minimal edge-cut (minimizing communication).

Consider Matrix-Vector Multiplication



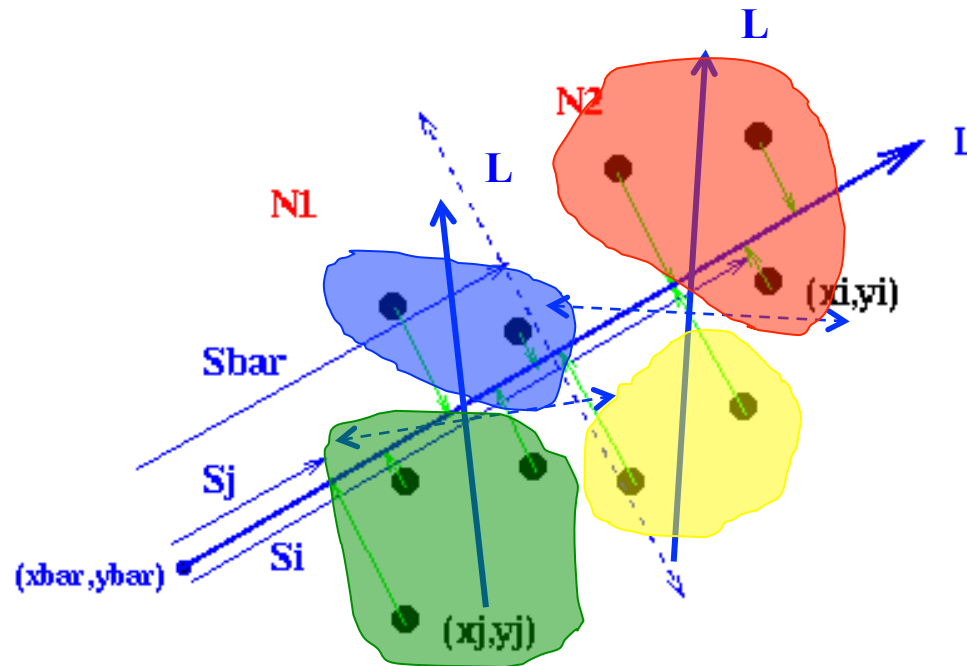
Is there a better partitioning than the straight forward with less edge-cut?





## 1. Recursive Inertial partitioning

Find the axis of minimal inertia and divide into two halves perpendicular to the axis.



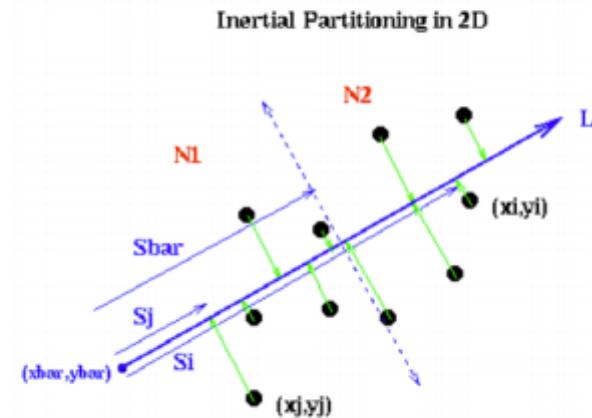
Orthogonal cut gives “best” edge-cut.  
Proceed in the same way with each half.



# 1. Recursive Inertial partitioning

- For a graph in 2D, choose line with half the nodes on one side and half on the other
  - In 3D, choose a plane, but consider 2D for simplicity
- Choose a line  $L$ , and then choose an  $L^\perp$  perpendicular to it, with half the nodes on either side

- 1)  $L$  given by  $a^*(x-xbar)+b^*(y-ybar)=0$ , with  $a^2+b^2=1$ ;  $(a,b)$  is unit vector  $\perp$  to  $L$
- 2) For each  $n_j = (x_j, y_j)$ , compute coordinate  $S_j = -b^*(x_j-xbar) + a^*(y_j-ybar)$  along  $L$
- 3) Let  $Sbar = \text{median}(S_1, \dots, S_n)$
- 4) Let nodes with  $S_j < Sbar$  be in  $N_1$ , rest in  $N_2$

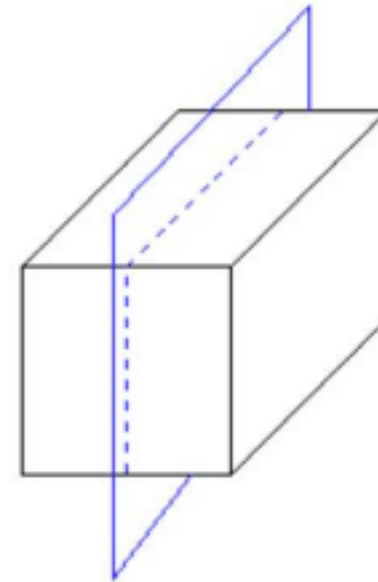


- Remains to choose  $L$



## 1. *Recursive Inertial partitioning*

- Mathematically, choose  $L$  to be a total least squares fit of the nodes
  - Minimize sum of squares of distances to  $L$
  - Equivalent to choosing  $L$  as axis of rotation that minimizes the moment of inertia of nodes





## 2. Recursive Spectral Bisection (RSB)

Assume that we have a graph with nodes  $V_I$  and edge weights  $W_{IJ}$  between nodes  $V_I$  and  $V_J$ . The graph can be represented with the Laplacian matrix  $L$ .

$$L_{IJ} = \begin{cases} -W_{IJ} & \text{if an edge } V_I \text{ to } V_J \\ \sum_K (W_{IK}) & \text{if } I=J \\ 0 & \text{otherwise} \end{cases}$$

## Theorem by Fiedler:

The second smallest eigenvalue to  $L$  satisfies

$$\lambda_2 = \min_{|x| \neq 0} \frac{\sum_{i,k} W_{ik} (x_i - x_k)^2}{\sum_i x_i^2}$$

and the minimum is attained for the corresponding eigenvector  $x = [x_1, x_2, \dots]$ .

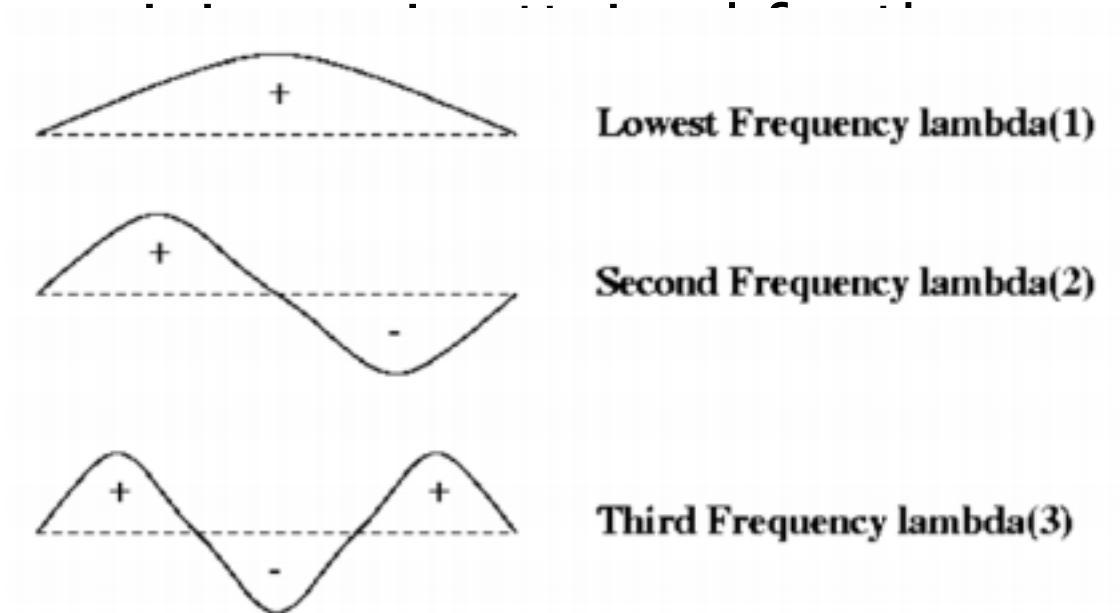
Note, to attain minimum  $x_i$  and  $x_k$  must be close if edge weight  $W_{ik}$  is large.

⇒ The second eigenvector gives geometrical information about the graph. Use this for partitioning of the graph.

## Theorem by Fiedler:

The second smallest eigenvalue to L satisfies

$$\lambda_2 = \min_{|x| \neq 0} \frac{\sum_{i,k} W_{ik} (x_i - x_k)^2}{\sum_i x_i^2}$$



partitioning of the graph.

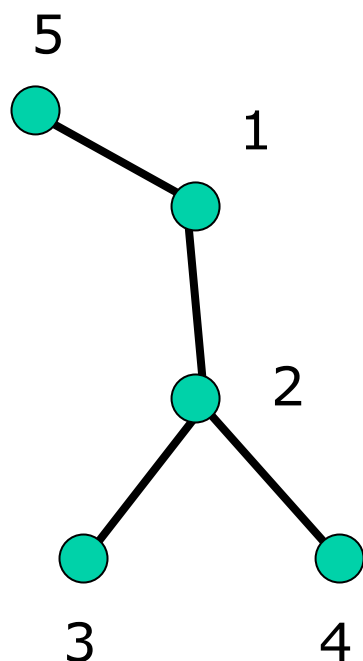


## Algorithm RSB:

1. Set up the Laplacian matrix.
2. Compute the eigenvector corresponding to the second smallest eigenvalue.
3. Sort the nodes according to the elements in eigenvector, then nodes with heavy edge weights will be close.
4. Divide the nodes in two halves according to the eigenvector.
5. Repeat 1-4 recursively with the two sets.

$$\lambda_2 = \min_{|x| \neq 0} \frac{\sum_{i,k} W_{ik} (x_i - x_k)^2}{\sum_i x_i^2}$$

Example:



Laplacian matrix =

2	-1	0	0	-1
-1	3	-1	-1	0
0	-1	1	0	0
0	-1	0	1	0
-1	0	0	0	1

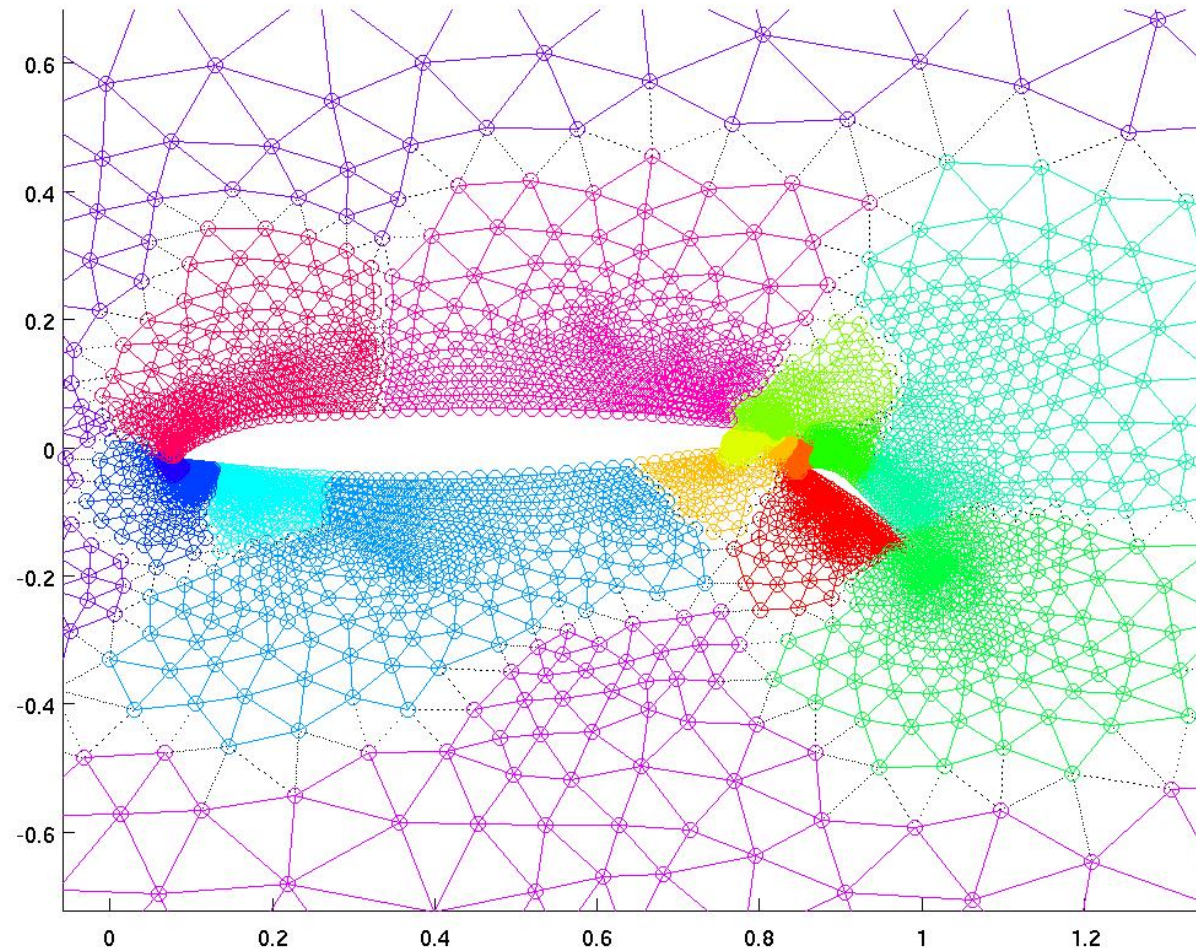
Eigenvalues =

1.0000	<b>0.5188</b>	0.0000	2.3111	4.1701
--------	---------------	--------	--------	--------

Eigenvectors =

0	<b>-0.3380</b>	-0.4472	0.7031	-0.4375
-0.0000	<b>0.2018</b>	-0.4472	0.3175	0.8115
0.7071	<b>0.4193</b>	-0.4472	-0.2422	-0.2560
-0.7071	<b>0.4193</b>	-0.4472	-0.2422	-0.2560
0.0000	<b>-0.7024</b>	-0.4472	-0.5362	0.1380

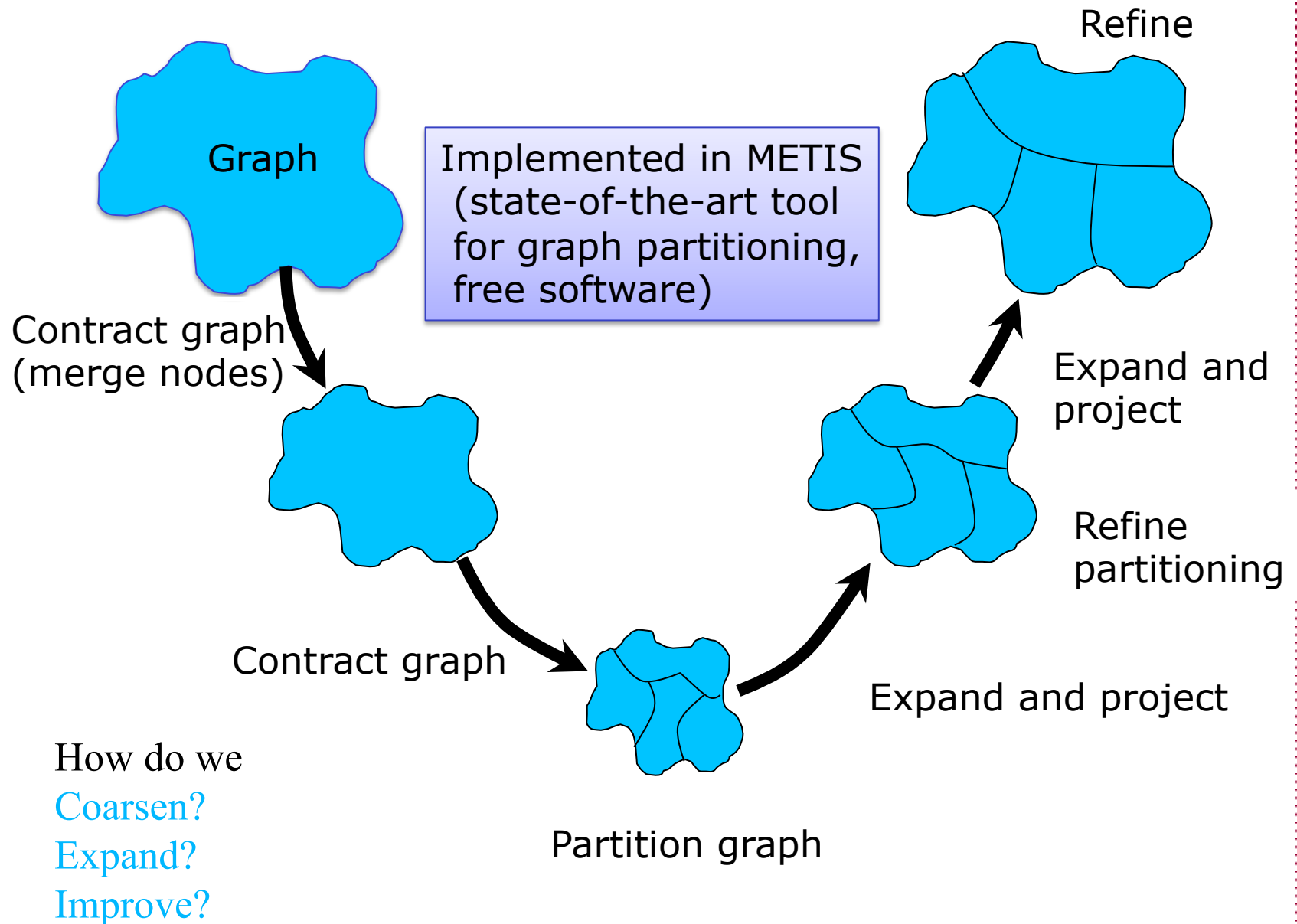




Example: Partitioning a FEM mesh  
(Note, the graph can be *very* large)

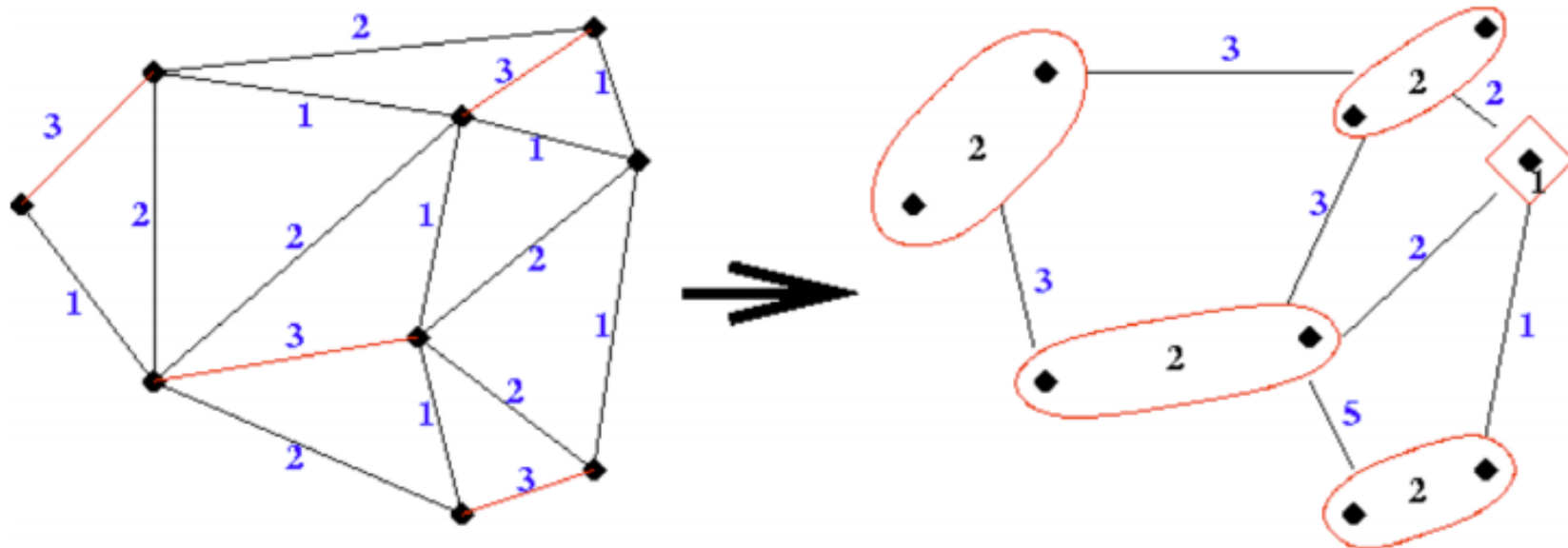


### 3. Multilevel partitioning



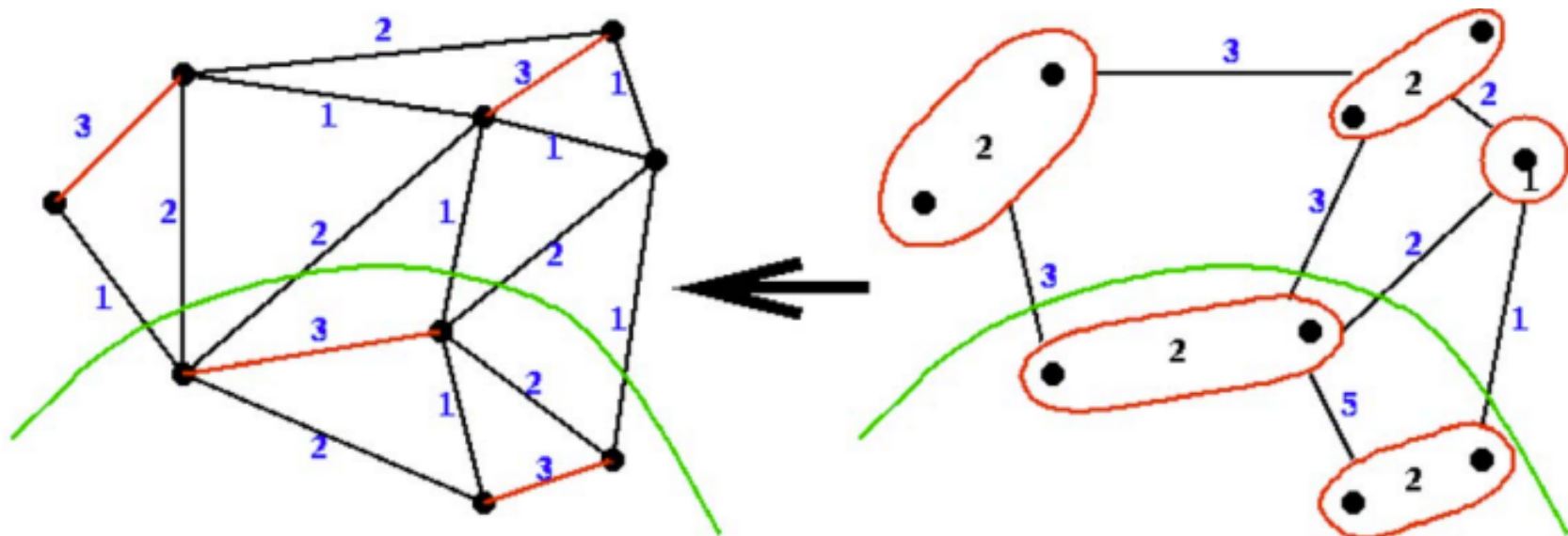
### 3. Multilevel partitioning--Coarsen

- Definition: A matching of a graph  $G(N,E)$  is a subset  $E_m$  of  $E$  such that no two edges in  $E_m$  share an endpoint
- Definition: A maximal matching of a graph  $G(N,E)$  is a matching  $E_m$  to which no more edges can be added and remain a matching



### 3. Multilevel partitioning—Partition

- Definition: An independent set of a graph  $G(N,E)$  is a subset  $N_i$  of  $N$  such that no two nodes in  $N_i$  are connected by an edge
- Definition: A maximal independent set of a graph  $G(N,E)$  is an independent set  $N_i$  to which no more nodes can be added and remain an independent set



### *3. Multilevel partitioning — Available Implementation*

- Multilevel Kernighan/Lin
  - METIS ([www.cs.umn.edu/~metis](http://www.cs.umn.edu/~metis))
  - ParMETIS - parallel version
- Multilevel Spectral Bisection
  - S. Barnard and H. Simon, “A fast multilevel implementation of recursive spectral bisection ...”, Proc. 6th SIAM Conf. On Parallel Processing, 1993
  - Chaco ([www.cs.sandia.gov/CRF/papers\\_chaco.html](http://www.cs.sandia.gov/CRF/papers_chaco.html))
- Hybrids possible
  - Ex: Using Kernighan/Lin to improve a partition from spectral bisection

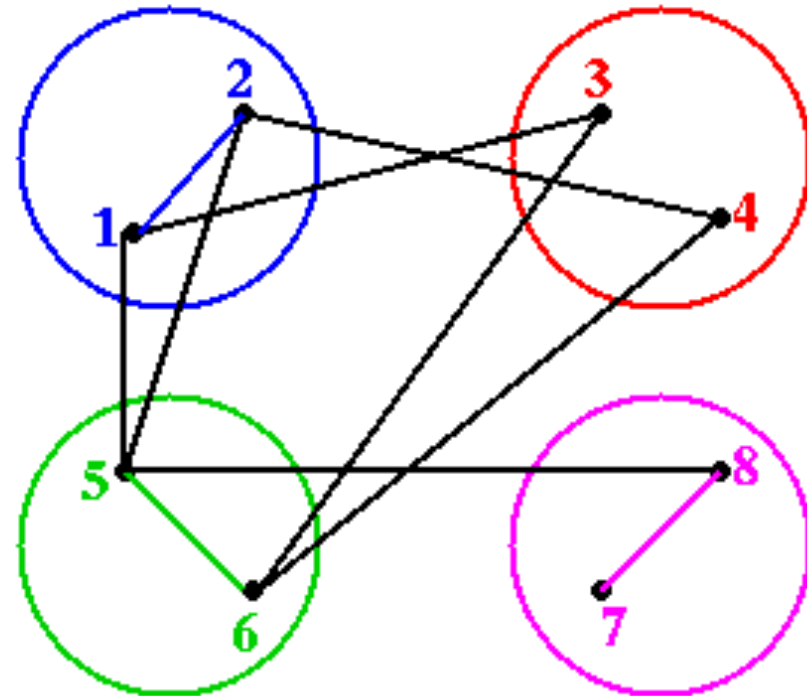
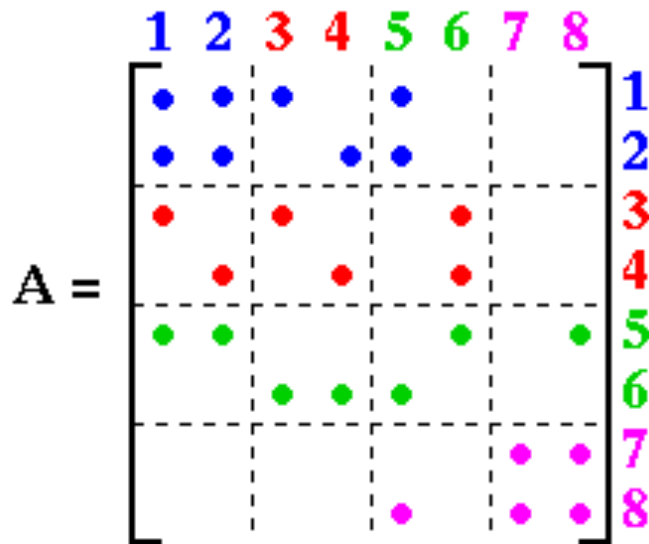
## 4. Bandwidth minimization

FEM  $\rightarrow$   $Ax=b$  where  $A$  is sparse  
( $<1\%$  non-zero) and very  
large ( $\sim 10^6$  rows).

$Ax=b$  is usually solved with an iterative method based on a sequence of Matrix-vector multiplications (e.g. Conjugate gradient method).

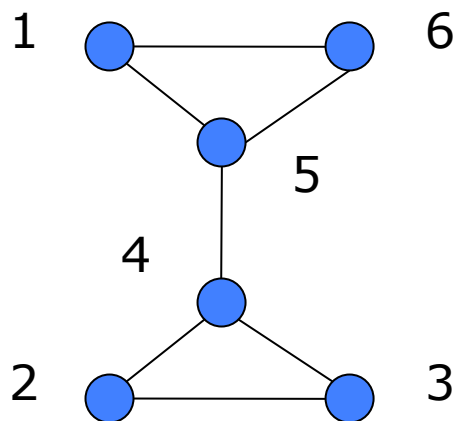
Need to parallelize the  $MxV$  row-wise. Block off-diagonal elements requires communication.





Minimize bandwidth => minimize  
communication and improve cache locality!

The bandwidth depends on the numbering of the nodes, consider the following graph and the corresponding matrix:

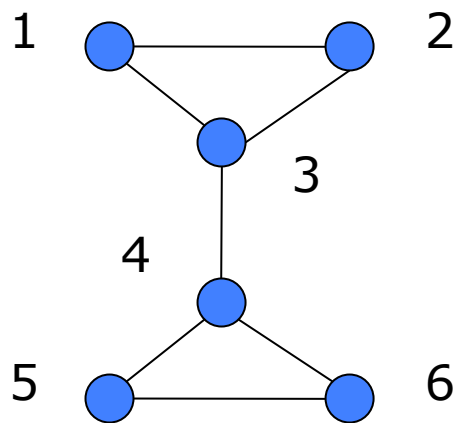


X				X	X
	X	X	X		
	X	X	X		
	X	X	X	X	
X			X	X	X
X				X	X

Bandwidth = 6



Then re-number the nodes:



X	X	X			
X	X	X			
X	X	X	X		
		X	X	X	X
			X	X	X
			X	X	X

Bandwidth = 4

## ***Reverse Cuthill-McKee algorithm:***

1. Find a root node, i.e., a node with a small number of edges => First node.
2. Number its neighbors, starting with fewest edges, then next fewest, etc.
3. Continue with lowest numbered node, with unnumbered neighbors, and number its neighbors as above.
4. When all nodes are numbered, reverse the ordering.

## ***Cuthill-McKee algorithm:***

First we choose a **peripheral vertex** (the vertex with the lowest **degree**)  $x$  and set  $R := (\{x\})$ .

Then for  $i = 1, 2, \dots$  we iterate the following steps while  $|R| < n$

- Construct the adjacency set  $A_i$  of  $R_i$  (with  $R_i$  the  $i$ -th component of  $R$ ) and exclude the vertices we already have in  $R$

$$A_i := \text{Adj}(R_i) \setminus R$$

- Sort  $A_i$  with ascending vertex order (**vertex degree**).
- Append  $A_i$  to the Result set  $R$ .

Implemented as `p=symrcm(A)` in Matlab, takes a matrix  $A$  and returns a permutation vector  $p$  minimizing the bandwidth of  $A(p,p)$ .

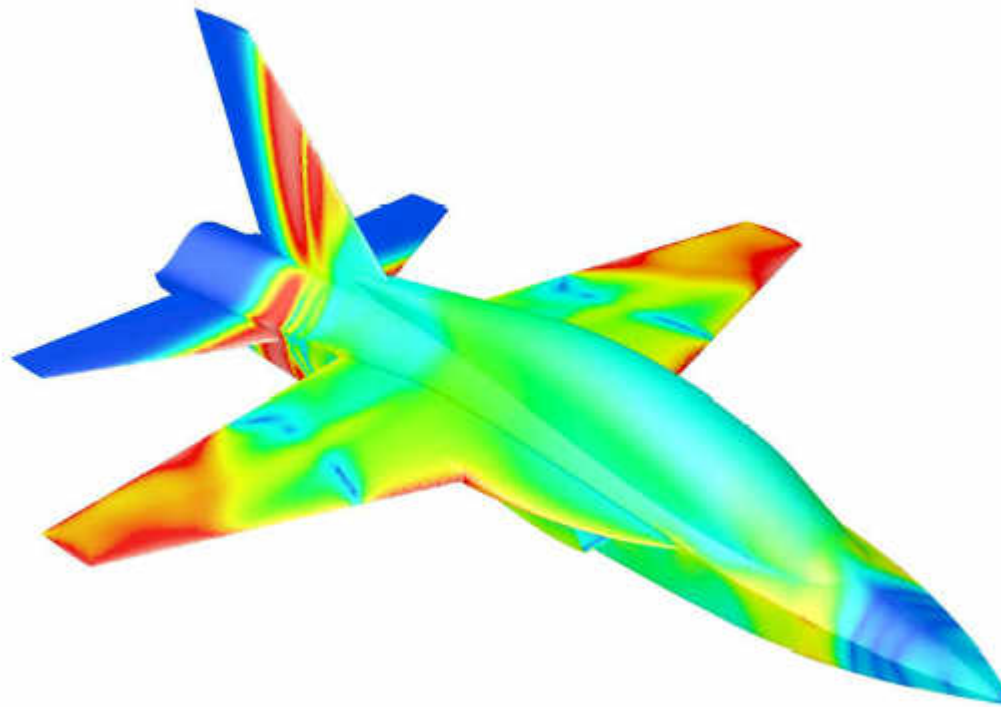
## ***Cuthill-McKee algorithm:***

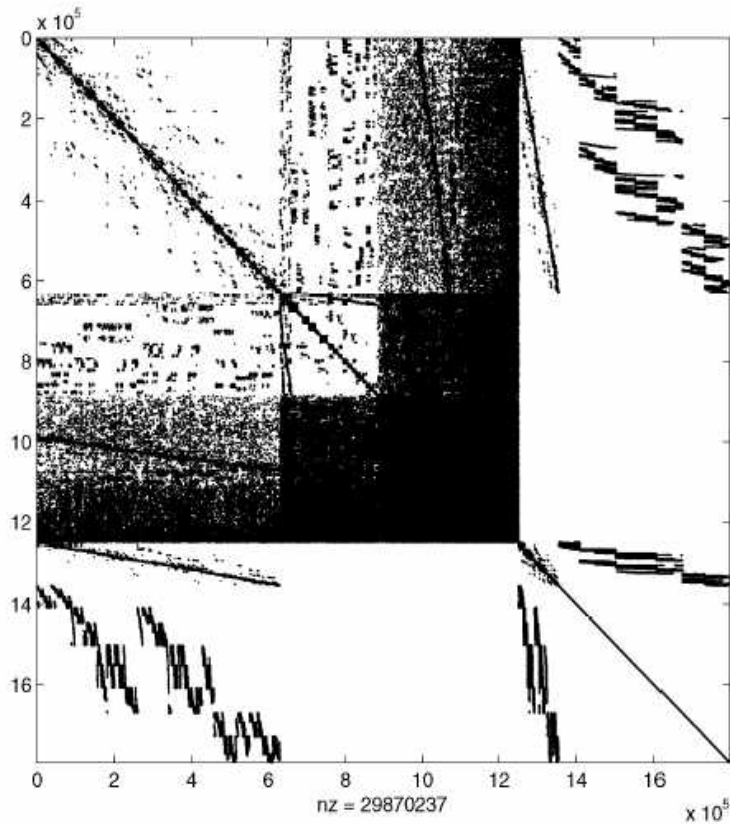
Read more on:

- [The Cuthill-McKee algorithm](#)
- [symrcm](#) in Matlab
- Reverse cuthill mckee in [scipy](#)
- [Matrix bandwidth and band matrix](#)

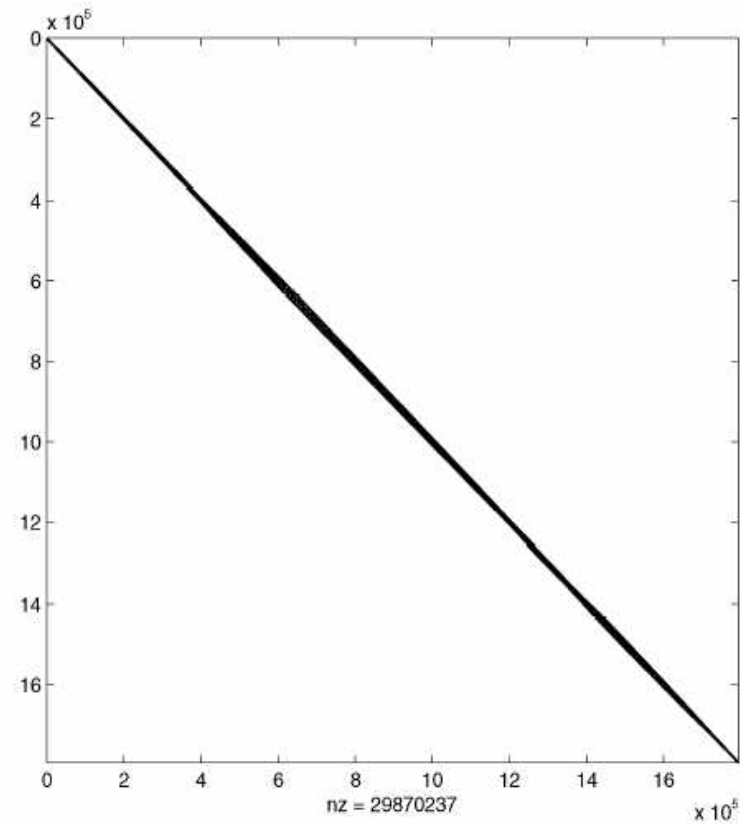
## Real applications: GEMS-project

Maxwell's equations discretized with FEM-grid around a fighter jet =>  $Ax=b$  with 1.8 million unknowns, solved with the CG method .





Original matrix



Bandwidth minimized

[ Ref: H. Löf, J. Rantakokko, *Algorithmic Optimization of a Conjugate Gradient Solver on Shared memory systems*, International Journal of Parallel, Emergent and Distributed Systems, Vol 21, 2006.]

# Ocean modeling, SMHI

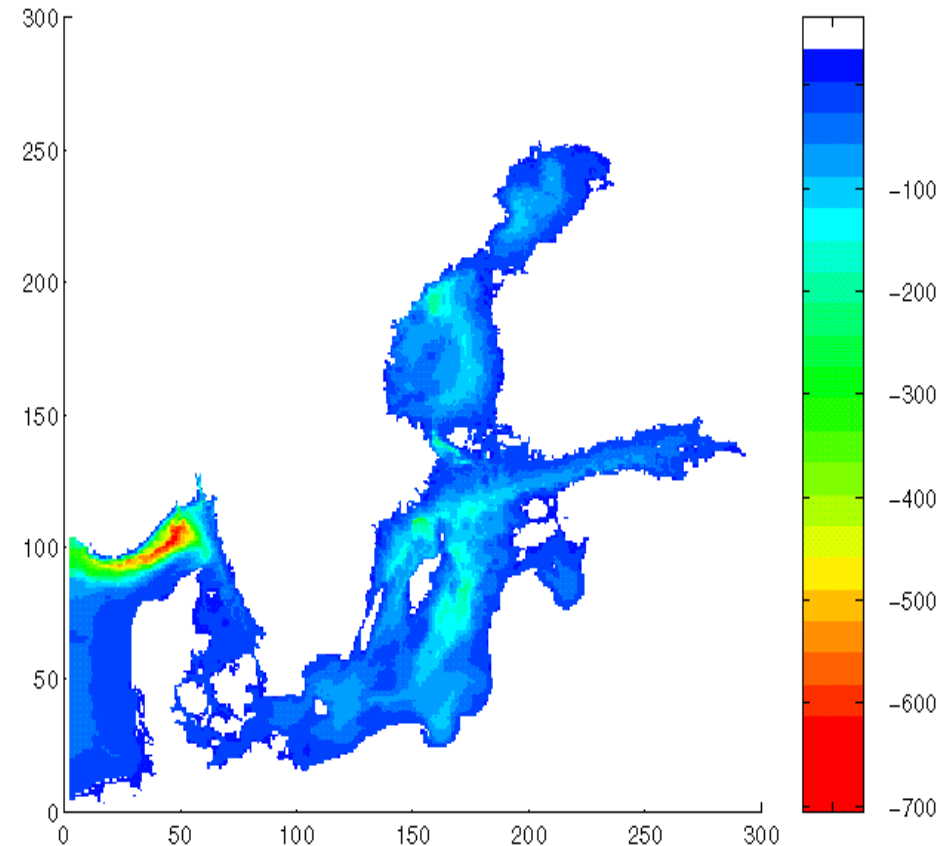
Baltic Sea, 300x300  
structured grid  
varying sea depth



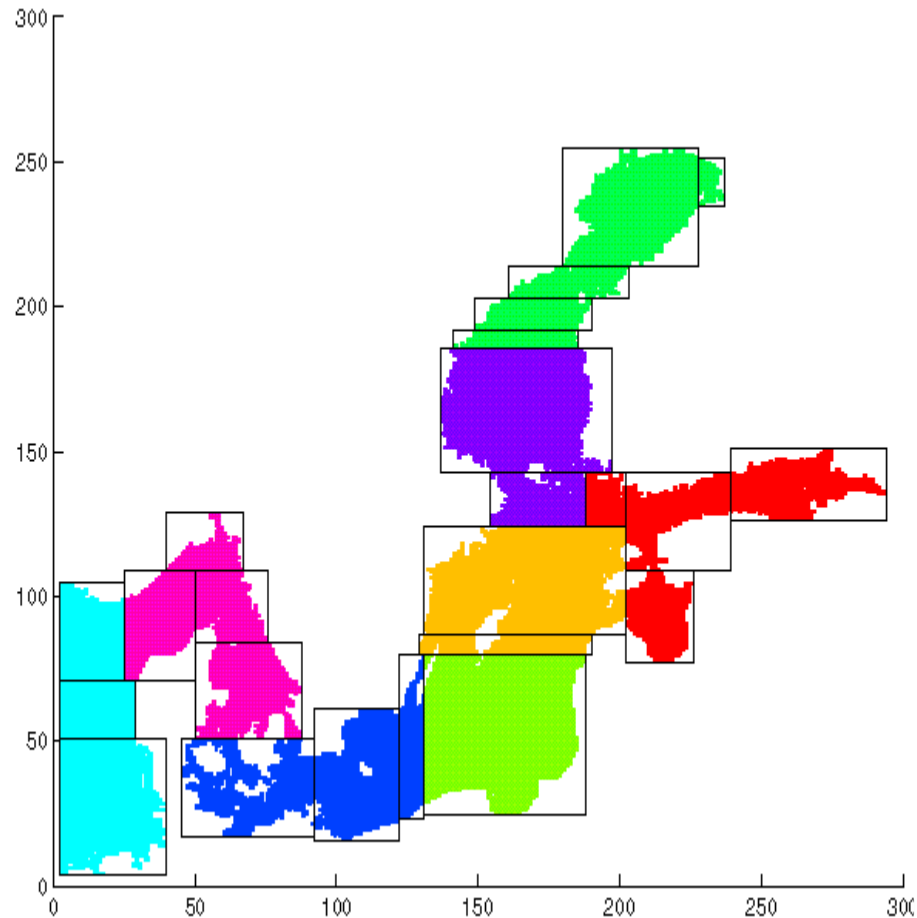
Need a block partition  
but have varying work  
load per grid point  
(zero for land points)

Minimize:

- load imbalance
- communication
- land points
- blocks



## Hybrid block-structured & graph partitioning method:



1. Cover the domain with  $m \times m$  blocks.
2. Remove blocks completely on land.
3. Shrink blocks to sea boundary.
4. Split blocks with large fraction of land and repeat 2-4.
5. Set up a graph for remaining blocks.
6. Compute Fiedler vector and sort.
7. Divide into two sets and split a block in division if necessary.
8. Proceed with steps 5-7 recursively.

[ Ref: T. Wilhelmsson, J. Schüle, J. Rantakokko, L. Funkquist, *Increasing Resolution and Forecast Length with a Parallel Ocean Model*, In: *Operational Oceanography, Implementation at the European and Regional Scales*, editor N.C. Fleming, Elsevier Oceanography Series, 2002.]