

# Solving linear system

**Francesco Sanfilippo**

Istituto Nazionale di Fisica Nucleare - Sezione Roma Tre



27 - 31 March 2017

## Introduction of the problem

- ① Direct solution of linear system  $Ax = b$
- ② Quadratic functional minimization

## Iterative solver

- ① Advantages
- ② Comparison of efficiency

## Checking the solution

- ① Limits, stability and efficiency of various algorithms
- ② Convergence criterions

## Accelerating the convergence

- ① Mixed precision algorithms
- ② Choosing a starting guess
- ③ Preconditioning the problem

## Solving similar problems at the same time

- ① Shifted problems  $A + \sigma \text{Id}$
- ② Deflating the problem

## Review of Parallelisation

- ① Distributed memory
- ② Shared memory
- ③ Vectors

## Gather/scatter approaches

→ 1+2 different examples of gathering of non-local data

## More specifically on parallelisation

- ① Communication/computation overlap
- ② Multithreading
- ③ Vectorization

## An example of a physical application

→ Lattice QCD

# Implicit and explicit residue

## Implicit residue

- The value of the residue can be implicitly computed looking  $|r_k|$
- This vector is automatically computed implicitly with:  $r_{k+1} = r_k - \alpha_k p_k$

## Explicit residue

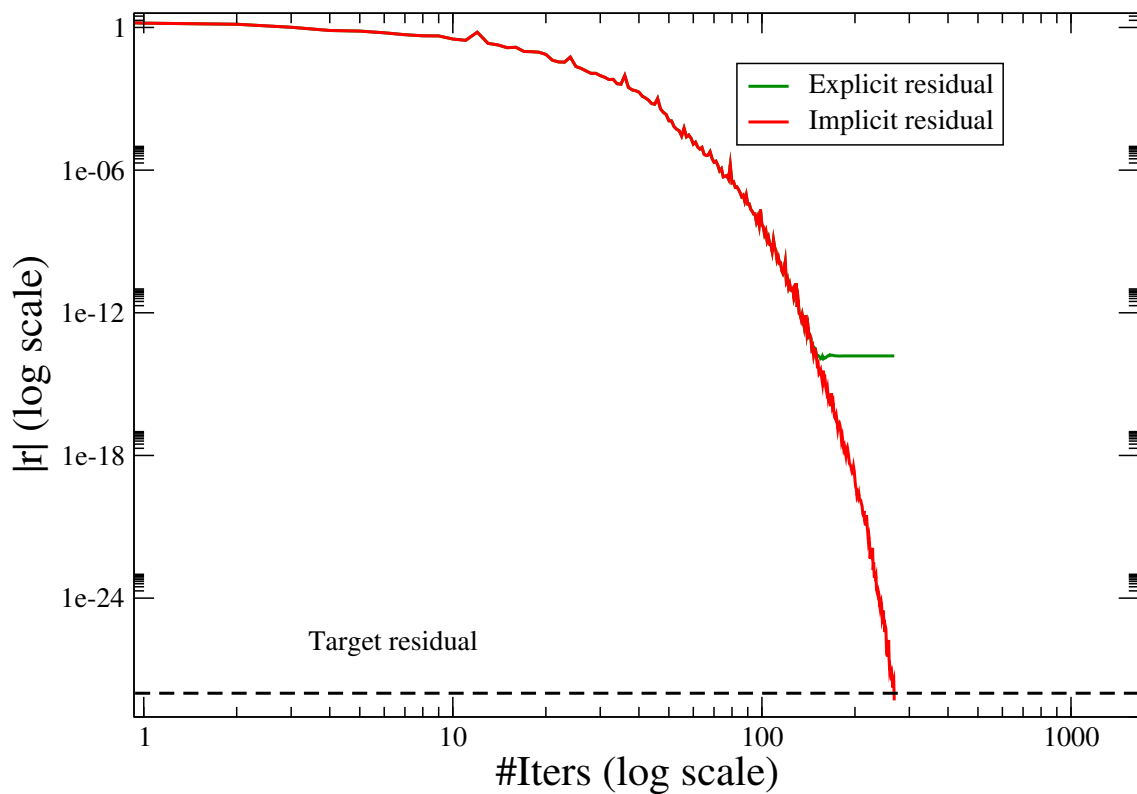
- One can also compute  $r_k^{expl} = Ax_k - b$  at the cost of an **additional application**
- This way is possible to compute  $F$  as well

## Numerical rounding

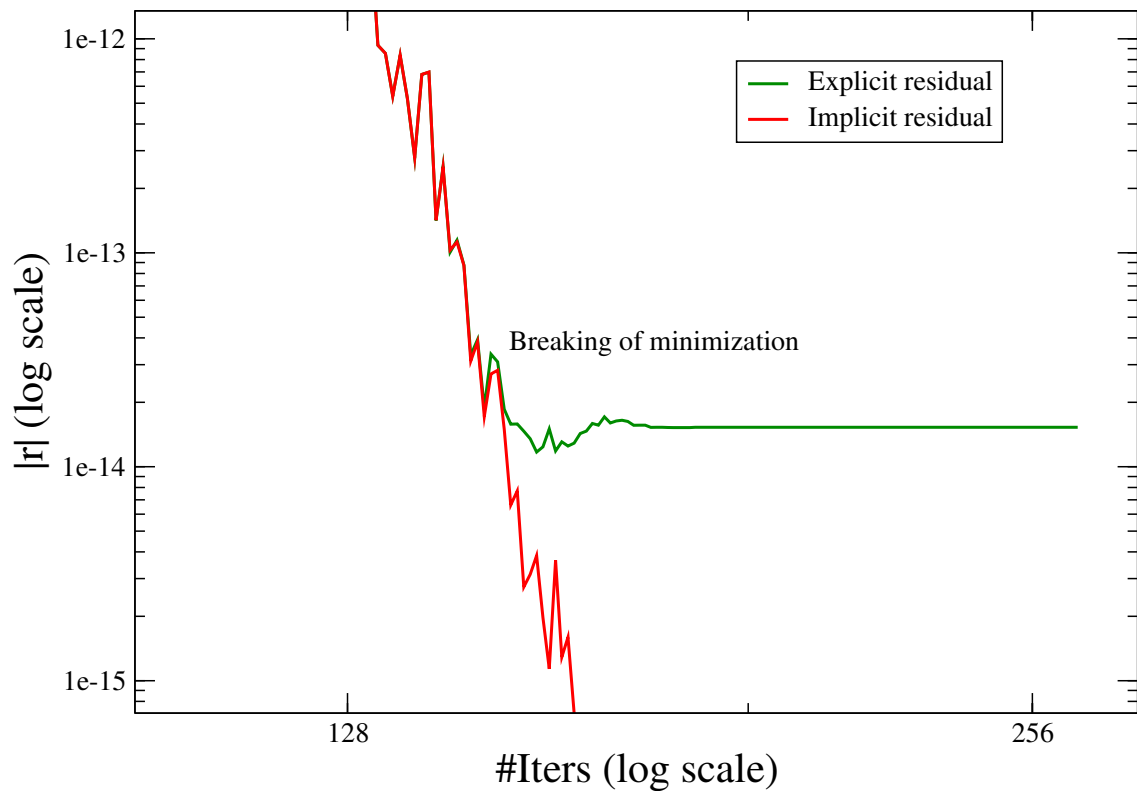
- As an effect of error accumulation,  $|r_k^{expl} - r_k^{impl}|$  will grow with  $k$
- This manifests the increase breaking of automatic conjugacy of generated  $p_k$
- After a certain number of iterations the whole mechanism of minimization of  $F$  breaks

**You cannot decrease  $|r|$  arbitrarily!**

Implicit and explicit residue in Conjugate gradient,  $|r|_{\text{targ}} = 10^{-28}$ ,  $n = 140$



Implicit and explicit residue in Conjugate gradient,  $|r|_{\text{targ}} = 10^{-28}$ ,  $N = 140$





## Improving the result beyond machine precision

Split the system  $\mathbf{A}x = b \rightarrow \mathbf{A}(x_0 + x_1 + \dots x_{n-1}) = b$

- First set  $b_0 = b$  and obtain an approximated solution  $x_0$  such that

$$\mathbf{A}x_0 = b_0 \quad (+r_0)$$

with a target relative residual  $|r_0| / |b_0| = \hat{r}_{targ}$  reachable by the conjugate gradient

- Then solve for  $b_1 = -r_0$ :

$$\mathbf{A}x_1 = b_1 \quad (+r_1)$$

where the relative residual  $|r_1| / |b_1| = |r_0| / |b_0|$  is reachable by the conjugate gradient

- Iterate until  $|r_{n-1}| / |b| = \hat{r}_{targ}$
- The solution is given by  $x = x_0 + x_1 + \dots + x_{n-1}$

### Problem

- When the residue  $|r| \ll |b|$  (high precision reached)  $r$  comes from a big cancellation
- If  $|r| / |b| \sim \epsilon$  (machine precision  $\sim 10^{-16}$ ),  $r_0^{expl}$  is badly computed and differs from true  $r_0$
- $x_1$  obtained from  $\mathbf{A}x_1 = b_1$  will not to improve  $x_0$ :

$$\mathbf{A}(x_0 + x_1) = b + r_0 - r_0^{expl} \neq b$$

# Mixed precision

## Higher precision

- Modern libraries offer support for emulated higher precision algebra
- e.g. gnu compiler `__float128` type implements quadruple precision (128 bit),  $\epsilon \sim 10^{-32}$

## Efficiency

- Emulation has a cost (order of magnitudes slower than hardware implementation)
- Running the whole solver in quadruple precision is costly and inefficient

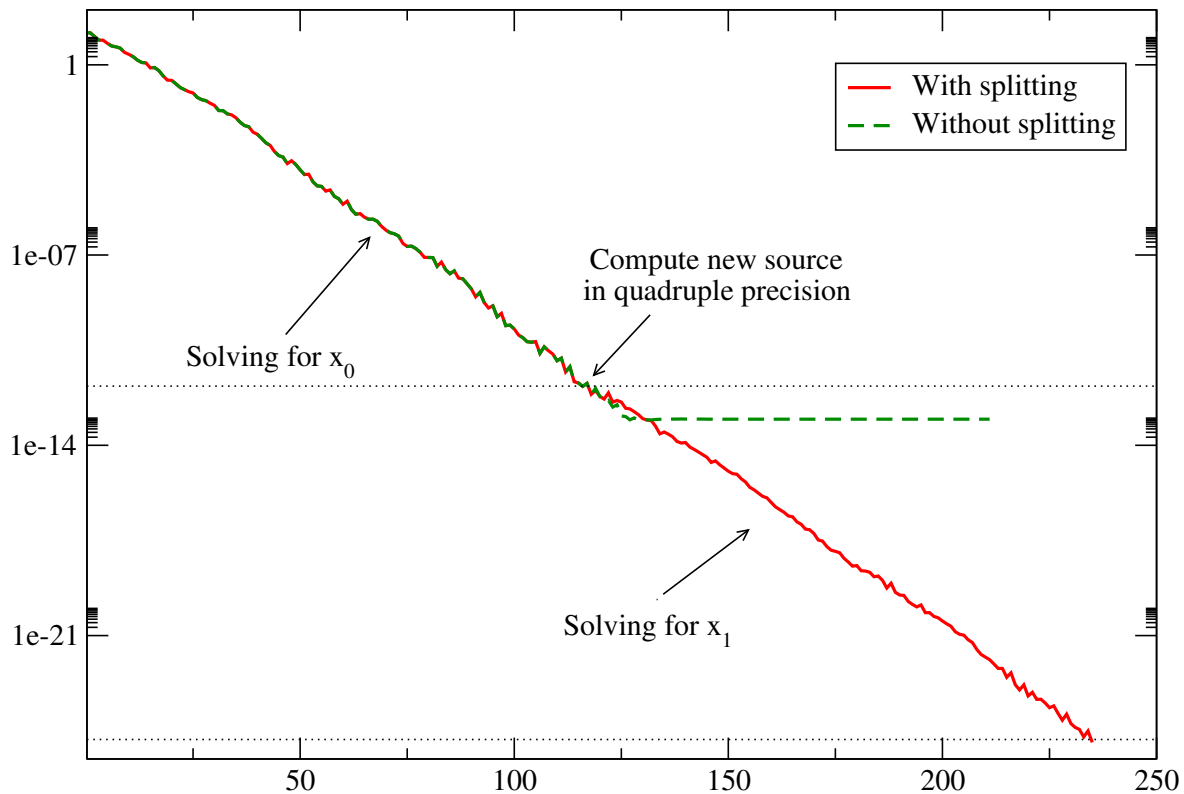
## Solution

- Problem comes from the fact that  $r_0^{expl} \neq r_0$ , so that

$$\mathbf{A}(x_0 + x_1) = b + r_0 - r_0^{expl} \neq b$$

- Only need to compute  $r_0^{expl}$  accurately, so that  $|r_0 - r_0^{expl}| \ll |r_0|$
- Quadruple precision is used only once in a while:
  - to compute  $b = -r_i$
  - to sum together  $x = x_0 + x_1 + \dots x_{n-1}$
- Ordinary precision is used through all the rest (inner solver)

# Mixed precision



## Acceleration via lower precision

The same game can be played the other way around:

- Use single (or even lower) precision in the inner solver
- Accumulate the external solution in higher precision

## Advantages

- Some architectures (most notably, GPU) lacks hardware support to double precision, or have much faster single precision
- Data needed to compute  $\mathbf{Ax}$  occupy less memory, so
  - more data coherence (optimize cache access)
  - faster loading from memory
- When data is split across different computing nodes, less communication is needed

## Disadvantages

- Need to store more information (lower and higher precision)
- At each restart Krylov space must be regenerated, can introduce unneeded overhead

## Initial guess

- So far we assumed to start from an initial guess  $x_0 = 0$  and iterate
- How can we incorporate a preliminary knowledge of an approximate solution  $x_{appr}$ , ?

## Splitting the problem again

Decompose the total solution

$$\mathbf{A}(x - x_{appr} + x_{appr}) = b$$

so that the problem can be rewritten as:

$$\mathbf{A}\Delta x = c$$

with:

$$c = b - r_{appr}, \quad r_{appr} = \mathbf{A}x_{appr}, \quad \Delta x = x - x_{appr}$$

## Recombine the problem

- After solving for  $\Delta x$  we sum  $x = x_{appr} + \Delta x$
- Equivalent to set from the beginning  $x_0 = x_{appr}$

## Which guess

For example we could have solved a similar system, such as:

$$\mathbf{A}x = b' \quad \text{or} \quad \mathbf{A}'x = b$$

# Non symmetric definite positive matrices

## Condition for convergence

- Minimization of the functional is guaranteed by the fact that  $F(x) = \frac{1}{2}x^T \mathbf{A}x - bx$  represents a paraboloid in a  $N$ -dimensional space, and  $x = \mathbf{A}^{-1}b$  locate the minimum
- Eigenvalues of  $\mathbf{A}$  must be real and positive ( $\mathbf{A}$  is symmetric and definite positive)
- What to do if  $\mathbf{A}$  does not satisfy the conditions?

## Conjugate Gradient Normal Equation (CGNE)

- For whatever  $\mathbf{A}$  (even rectangular) the matrix  $\mathbf{B} = \mathbf{A}^t \mathbf{A}$  is symmetric and definite positive, in facts eigenvalues of  $\mathbf{B}$  are the square modulo of those of  $\mathbf{A}$
- Equation obtained multiplying both sides of  $\mathbf{A}x = b$  by  $\mathbf{A}^t$  has the same solution  $x$ :

$$\mathbf{A}^t \mathbf{A}x = \mathbf{A}^t b$$

- Get  $x$  by solving  $\mathbf{B}$  for the vector  $c = \mathbf{A}^t b$ , as:  $\mathbf{B}x = c$

## Disadvantages

- Matrix  $\mathbf{B} = \mathbf{A}^t \mathbf{A}$  condition number is the square of that of  $\mathbf{A}$
- Applying  $\mathbf{B}$  cost: twice of applying  $\mathbf{A}$  to  $x$

## Variation: biconjugate (stabilized) gradient algorithm

### Problem

- When  $\mathbf{A}$  is not symmetric & positive definite  $\mathbf{A}$  is not a quadratic form
- (when the space is complex, symmetric  $\rightarrow$  hermitian)
- Minimization breaks down

### Biconjugate stabilized

- Instead of minimizing  $F$  we can find other functional forms
- Biconjugate stabilized algorithm minimizes also non-symmetric definite positive matrix
- Pro:
  - Works for generic matrix
  - Convergence typically requires less iterations
- Con:
  - Two applications of  $\mathbf{A}$  per iteration needed
  - Is more turbulent (residue can jump)
  - Convergence is not guaranteed in all the cases

### Other options

- Minimum residual method (MINRES), Generalized minimal residual method (GMRES) works to minimize  $|r_k|$  at each iteration
- Symmetric LQ method (SYMMLQ)

Convergence can be faster, but typically is never guaranteed for generic problem

# Preconditioning the problem

## Preconditioner

Introduce an operator  $\mathbf{P}$  such that:

- $\kappa(\mathbf{AP}^{-1}) < \kappa(\mathbf{A})$  the product  $\mathbf{AP}^{-1}$  has a lower condition number of  $\mathbf{A}$
- computing  $\mathbf{AP}^{-1}\mathbf{x}$  is not much more expensive than computing  $\mathbf{Ax}$

## Preconditioned equation

$$\mathbf{AP}^{-1}\mathbf{Px} = b$$

- Defining  $y = \mathbf{Px}$  solve:  $\mathbf{AP}^{-1}y = b$
- Then solve:  $\mathbf{Px} = y$  and obtain  $x$

## Definite-positivity

- If  $\mathbf{P}^{-1}$  does not commute with  $\mathbf{A}$ ,  $\mathbf{AP}^{-1}$  is not symmetric (even if  $\mathbf{A}$  is)
- Solve the problem for  $\mathbf{C} = \mathbf{P}^{-1}\mathbf{AP}^{-1}$  and  $d = \mathbf{P}^{-1}b$ :

$$\mathbf{C}y = d$$

Alternative preconditioned equation: just multiply the matrix from the left

$$\mathbf{P}^{-1}\mathbf{Ax} = \mathbf{P}^{-1}b = c$$



# Which preconditioner?

## Desiderata

- System for  $\mathbf{AP}^{-1}$  (or  $\mathbf{P}^{-1}\mathbf{A}$ ) must be simpler to be solved than  $\mathbf{A}$
- $\mathbf{P}$  should therefore approximate  $\mathbf{A}$  in some sense
- Computing  $\mathbf{P}^{-1}\mathbf{v}$  must be cheap

## Generic possibilities

- Jacobi preconditioner  $\mathbf{P} = \text{diag}(\mathbf{A})$   
good if the system is dominated by its diagonal
- Incomplete  $LU$  factorization: we force matrix  $\mathbf{U}$  to be zero where  $\mathbf{A}$  was good for sparse matrices

## More specific possibilities for sparse matrices

- Consider a simplified version of the matrix retaining its “physical” features
- Domain decomposition: consider separately sub-blocks of  $\mathbf{A}$   
(of size possibly related to a physical property of the system)

# Focus on a particular problem - Laplace equation in $D$ dimensions

## In the continuum

Setting  $\sigma = \text{const}$ ,  $\Delta = \sum_{\mu} \nabla_{\mu}^2$ :

$$(\sigma - \Delta) f(x) = b(x)$$

## Discretization

$$\begin{cases} f(x) & \rightarrow f_i \\ b(x) & \rightarrow b_i \\ (\nabla_{\mu}^2 f)(x) & \rightarrow \frac{1}{2} (f_{i+\hat{\mu}} + f_{i-\hat{\mu}} - 2f_i) \end{cases}$$

## Discrete problem

$$\mathbf{M}_{ij} f_j = b_i$$

Reduced to a linear system for matrix  $\mathbf{M}$ :

$$\mathbf{M}_{ij} = (\sigma + D) \delta_{ij} - \frac{1}{2} \sum_{\mu} (\delta_{i,j+\hat{\mu}} + \delta_{i,j-\hat{\mu}})$$

# Laplace equation in $D$ dimensions: exact solution

## In $x$ space

$$\mathbf{M}_{ij} f_j = b_i, \quad \mathbf{M}_{ij} = (\sigma + D) \delta_{ij} - \frac{1}{2} \sum_{\mu} (\delta_{i,j+\hat{\mu}} + \delta_{i,j-\hat{\mu}})$$

## Passing to Fourier space

$$\begin{cases} f_j = \frac{1}{V} \sum_k e^{ijk} f(k) \\ \mathbf{M}_{jl} = \frac{1}{V} \sum_k e^{i(j-l)k} \left( \sigma + D - \sum_{\mu} \cos k_{\mu} \right) \end{cases}$$

## Exact solution

$$f_j = \sum_k e^{-ijk} \mathbf{M}^{-1}(k) b(k) = \sum_k \frac{e^{-ijk} b(k)}{\sigma + D - \sum_{\mu} \cos k_{\mu}}$$

Good to check approximate solution!

## Eigenvalues, eigenvectors

- $M_{ij} = (\sigma + D) \delta_{ij} - \frac{1}{2} \sum_{\mu} (\delta_{i,j+\hat{\mu}} + \delta_{i,j-\hat{\mu}})$  is symmetric
- Due to translation invariance, plane waves are eigenvectors
- Eigenvalues labelled by corresponding momenta:

$$\lambda(k) = \sigma + D - \sum_{\mu} \cos k_{\mu}$$

- If  $\sigma = 0$  the system has a zero eigenvalue

## Matricially

For  $N = 6$ ,  $D = 1$  ( $d = \sigma + D$ ,  $s = -1/2$ )

$$\begin{pmatrix} d & s & 0 & 0 & 0 & s \\ s & d & s & 0 & 0 & 0 \\ 0 & s & d & s & 0 & 0 \\ 0 & 0 & s & d & s & 0 \\ 0 & 0 & 0 & s & d & s \\ s & 0 & 0 & 0 & s & d \end{pmatrix}$$

# Decomposing the problem - Even Odd precondition

## Splitting parity coupling (red-black)

$$\mathbf{M} = \mathbf{M}^{ee} + \mathbf{M}^{oe} + \mathbf{M}^{eo} + \mathbf{M}^{oo}$$

where e.g.  $\mathbf{M}^{eo}$  is different from 0 only between even and odd sites:

$$\mathbf{M}_{i,j}^{eo} = \mathbf{M}_{i,j} \delta_{p(i),e} \delta_{p(j),o}, \quad p(i) = \left( \sum_{\mu} x_i \right) \bmod 2, \quad e = 0, o = 1$$

and similarly

$$\mathbf{v} = \mathbf{v}^e + \mathbf{v}^o$$

## Rewriting the system

In this way the system is rewritten as 
$$\begin{cases} \mathbf{M}^{ee} \mathbf{x}^e + \mathbf{M}^{eo} \mathbf{x}^o = \mathbf{b}^e \\ \mathbf{M}^{oe} \mathbf{x}^e + \mathbf{M}^{oo} \mathbf{x}^o = \mathbf{b}^o \end{cases}$$

## Decoupling even solution

- Some algebra shows: 
$$\begin{cases} (\mathbf{M}^{ee} - \mathbf{M}^{eo} \frac{1}{\mathbf{M}^{oo}} \mathbf{M}^{oe}) \mathbf{x}^e = \mathbf{b}^e - \mathbf{M}^{eo} \frac{1}{\mathbf{M}^{oo}} \mathbf{b}^o \\ \mathbf{x}^o = \frac{1}{\mathbf{M}^{oo}} (\mathbf{b}^o - \mathbf{M}^{oe} \mathbf{x}^e) \end{cases}$$
- The equation for  $\mathbf{x}^e$  decouples from that of  $\mathbf{x}^o$

## Advantages of solving $x^e$

$$\left( M^{ee} - M^{eo} \frac{1}{M^{oo}} M^{oe} \right) x^e = b^e - M^{eo} \frac{1}{M^{oo}} b^o$$

### Cost

Application of  $M^{eo} \frac{1}{M^{oo}} M^{oe}$  and  $M^{ee}$  has the same cost of  $M$  but:

- $M_{ij}^{oo} = \delta_{ij}d$  so  $\frac{1}{M_{ij}^{oo}} = \frac{1}{d} M_{ij}$
- half of the data is necessary (less communications, more data coherence)
- the condition number has decreased

### Further preconditioning

- Additional optimization can be achieved factorizing  $M^{ee}$

$$\left( 1 - \frac{1}{M^{ee}} M^{eo} \frac{1}{M^{oo}} M^{oe} \right) x^e = \frac{1}{M^{ee}} \left( b^e - M^{eo} \frac{1}{M^{oo}} b^o \right)$$

### Solver

The system is still symmetric, so conjugate gradient solver usable

## Reordering

Putting first all even sites and then odd,

$$M = \begin{pmatrix} M^{ee} & M^{oe} \\ M^{eo} & M^{oo} \end{pmatrix}, \quad v = \begin{pmatrix} v_e \\ v_o \end{pmatrix}$$

## Good for performance

Not necessary to perform this reordering, but improve performance

$$\begin{pmatrix} d & s & 0 & 0 & 0 & s \\ s & d & s & 0 & 0 & 0 \\ 0 & s & d & s & 0 & 0 \\ 0 & 0 & s & d & s & 0 \\ 0 & 0 & 0 & s & d & s \\ s & 0 & 0 & 0 & s & d \end{pmatrix} \rightarrow \begin{pmatrix} d & 0 & 0 & 0 & s & s \\ 0 & d & 0 & s & 0 & s \\ 0 & 0 & d & s & s & 0 \\ 0 & s & s & d & 0 & 0 \\ s & 0 & s & 0 & d & 0 \\ s & s & 0 & 0 & 0 & d \end{pmatrix}$$

Access to a more compact piece of memory

# Domain decomposition

- What we saw is just an example of a Schur decomposition on the  $2 \times 2$  blocked matrix
- Other decompositions can come to your mind
- In our case for example

$$\left( \begin{array}{ccc|ccc} d & s & 0 & 0 & 0 & s \\ s & d & s & 0 & 0 & 0 \\ \hline 0 & s & d & s & 0 & 0 \\ 0 & 0 & s & d & s & 0 \\ 0 & 0 & 0 & s & d & s \\ s & 0 & 0 & 0 & s & d \end{array} \right) = \begin{pmatrix} D_{11} & \Omega_{12} \\ \Omega_{21} & D_{22} \end{pmatrix}$$

- In this case

$$\left( D_{11} - \Omega_{12} \frac{1}{D_{22}} \Omega_{21} \right) x^1 = b^1 - \Omega_{12} \frac{1}{D_{22}} b^2$$

- Computing  $t = \frac{1}{D_{22}} v$  requires actually to solve the  $D_{22}t = v$  (micro-system) but
  - the system is much smaller
  - in a parallel program, no communication between different domains needed
- Can be seen as “gluing together” the solution obtained on subdomain, and correcting it
  - Pro if sub-blocks size correctly chosen, few macro iteration needed
  - Con have to know the scale at which the micro-solution approximate well the total one



# Solving multiple similar systems

## Naive solution

Use the solution of a problem as first guess for the next one

## Shifted solver

If the problem is of the kind

$$(\mathbf{A} + \sigma)x = b$$

multi-shifted solver exists, e.g. M-CG, that can solve all the system for  $\sigma_1, \sigma_2, \dots$  simultaneously with very little overhead

## Deflation

If we are interested in multiple source  $b_1, b_2, \dots$  use algorithm (EigCG...) that

- while solving  $\mathbf{A}x = b$  finds the lowest eigenvector of the system
- for each sources  $b_1, b_2, \dots$  one eigenvector is found
- and can be removed by the spectra
- accelerating following solution