

4-9-2021

Lecture 09: Hierarchically Low Rank and Kronecker Methods

Rio Yokota
Tokyo Institute of Technology, rioyokota@gsic.titech.ac.jp

Follow this and additional works at: <https://scholarworks.uark.edu/mascsls>

 Part of the [Algebra Commons](#), [Control Theory Commons](#), [Non-linear Dynamics Commons](#), [Numerical Analysis and Computation Commons](#), and the [Numerical Analysis and Scientific Computing Commons](#)

Citation

Yokota, R. (2021). Lecture 09: Hierarchically Low Rank and Kronecker Methods. *Mathematical Sciences Spring Lecture Series*. Retrieved from <https://scholarworks.uark.edu/mascsls/16>

This Video is brought to you for free and open access by the Mathematical Sciences at ScholarWorks@UARK. It has been accepted for inclusion in Mathematical Sciences Spring Lecture Series by an authorized administrator of ScholarWorks@UARK. For more information, please contact scholar@uark.edu.

University of Arkansas Department of Mathematical Sciences
46th Spring Lecture Series

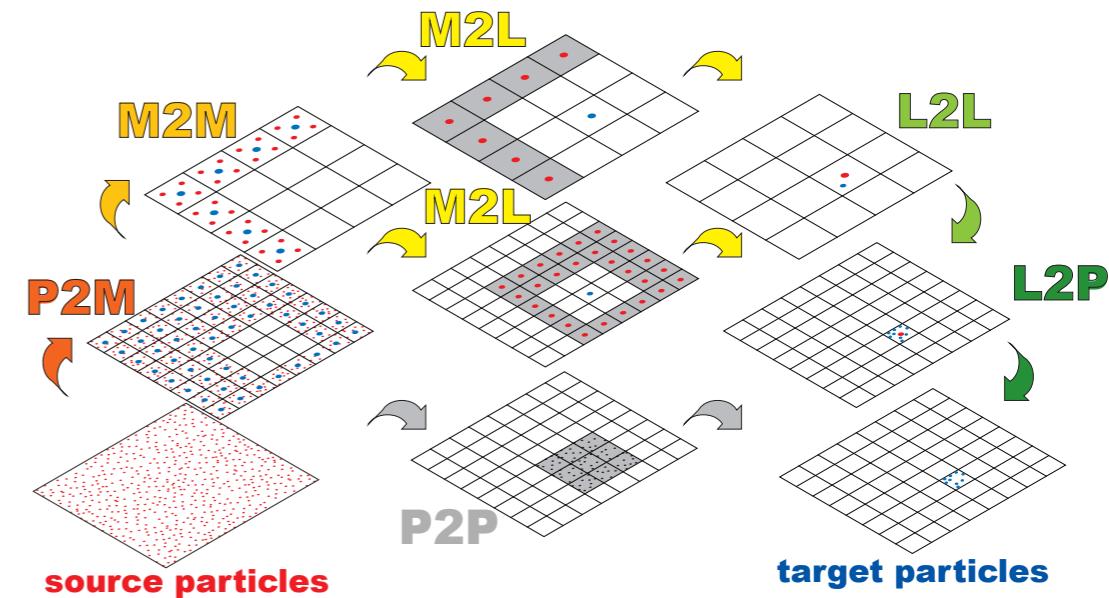
Hierarchically Low Rank and Kronecker Methods

Tokyo Institute of Technology
Rio Yokota
rioyokota@gsic.titech.ac.jp

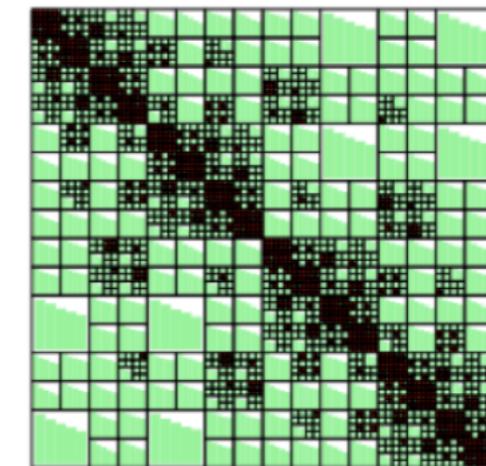


What I will be talking about today

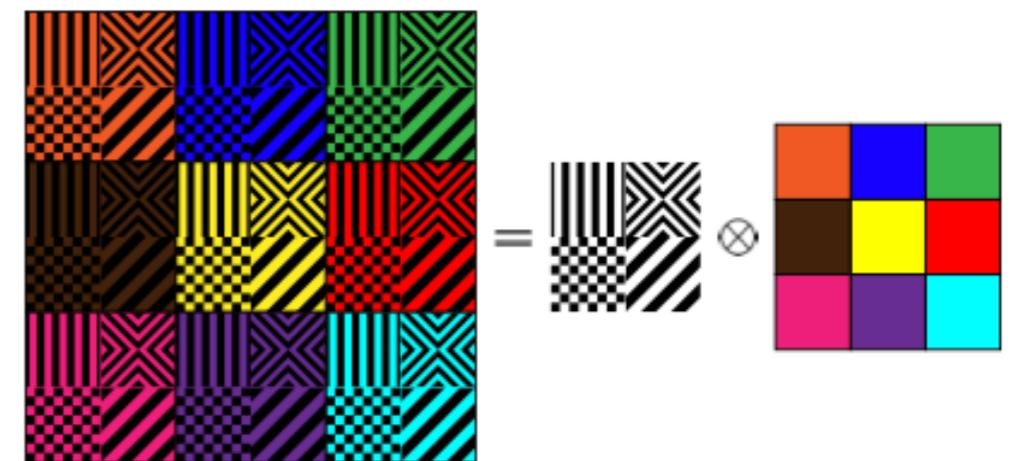
I. Fast multipole methods



2. Hierarchical low-rank matrices



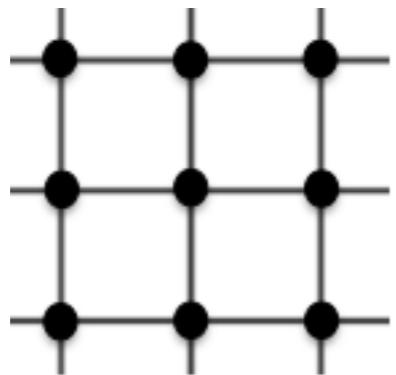
3. Kronecker factorization



Fast Multipole Methods

Structure of matrices

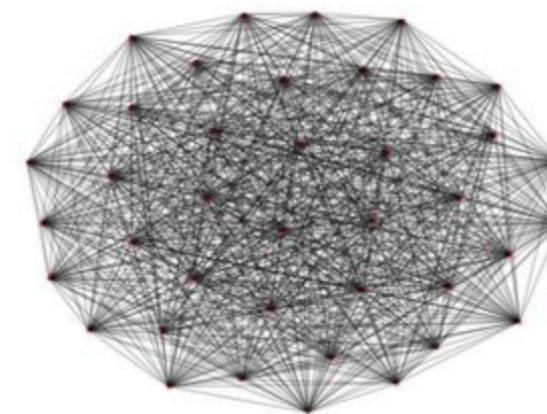
Sparse



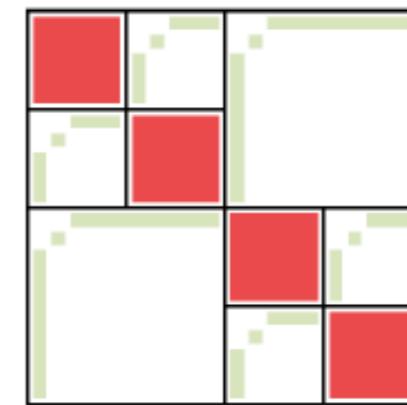
locally connected



Dense



fully connected



grouping based on connectivity

grouping based on proximity

Hierarchical N-body methods

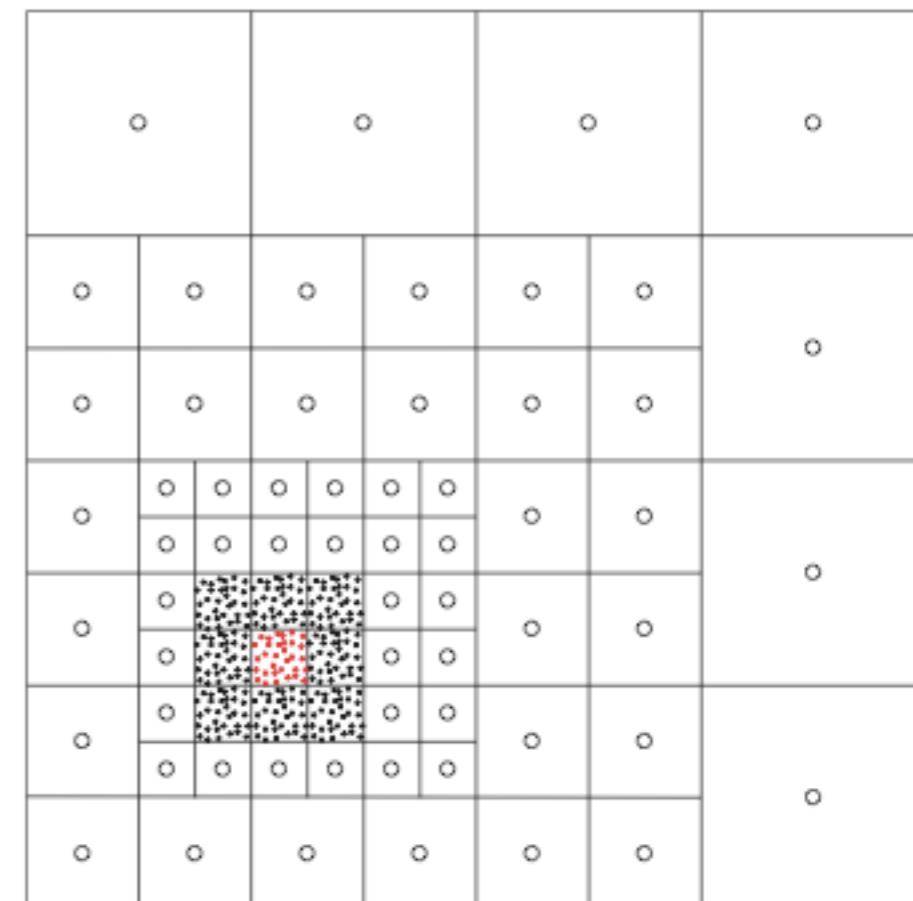
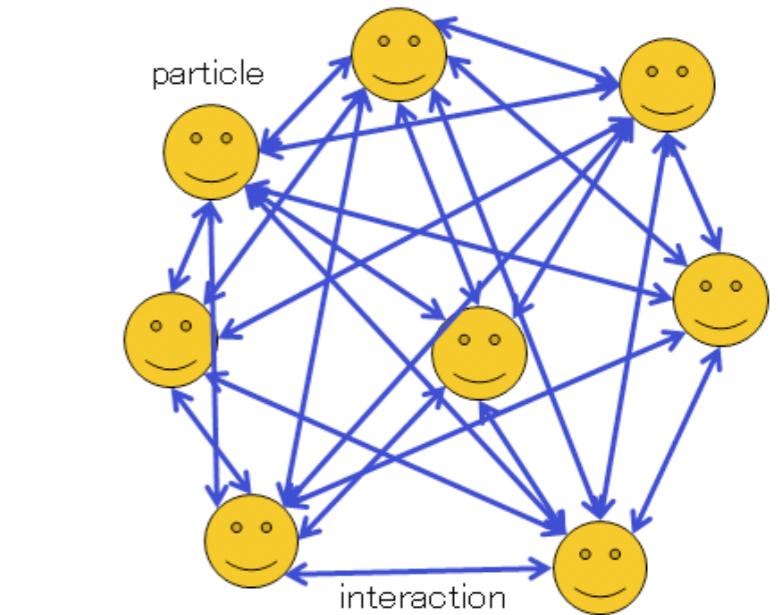
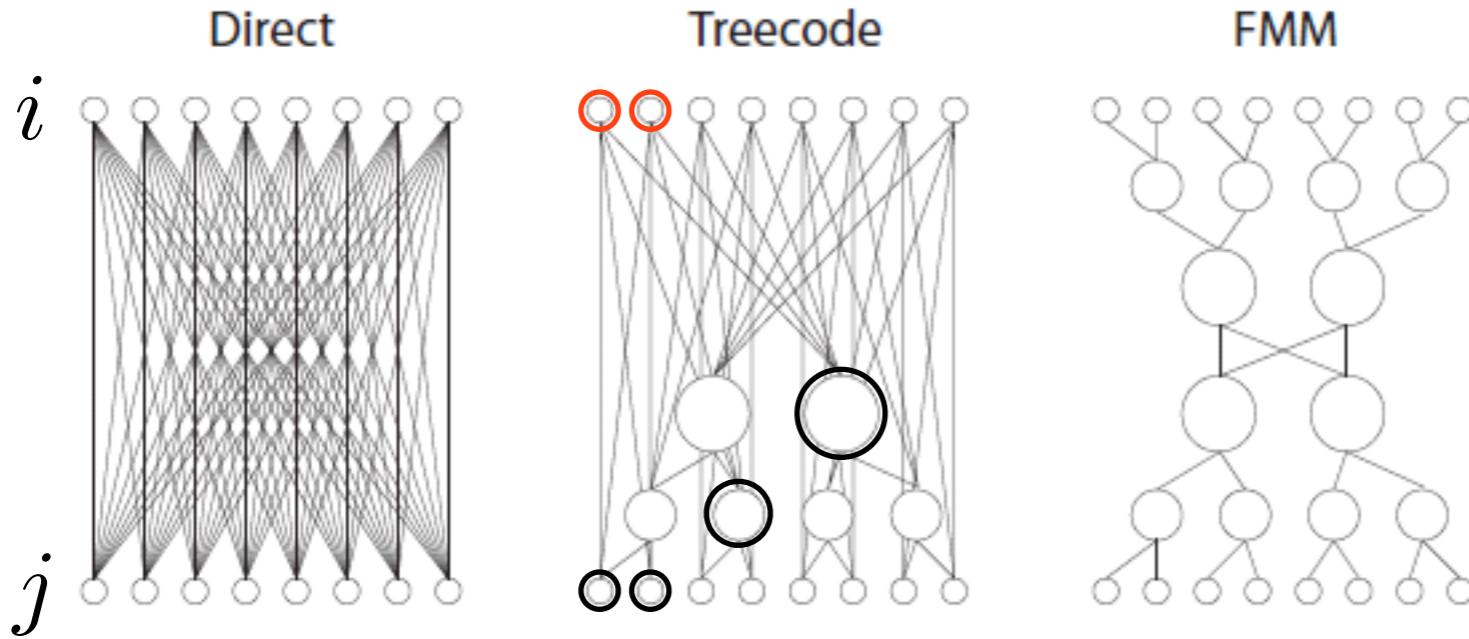
Particles interact with each other
Stars, Galaxies, Atoms, etc.

Computational cost

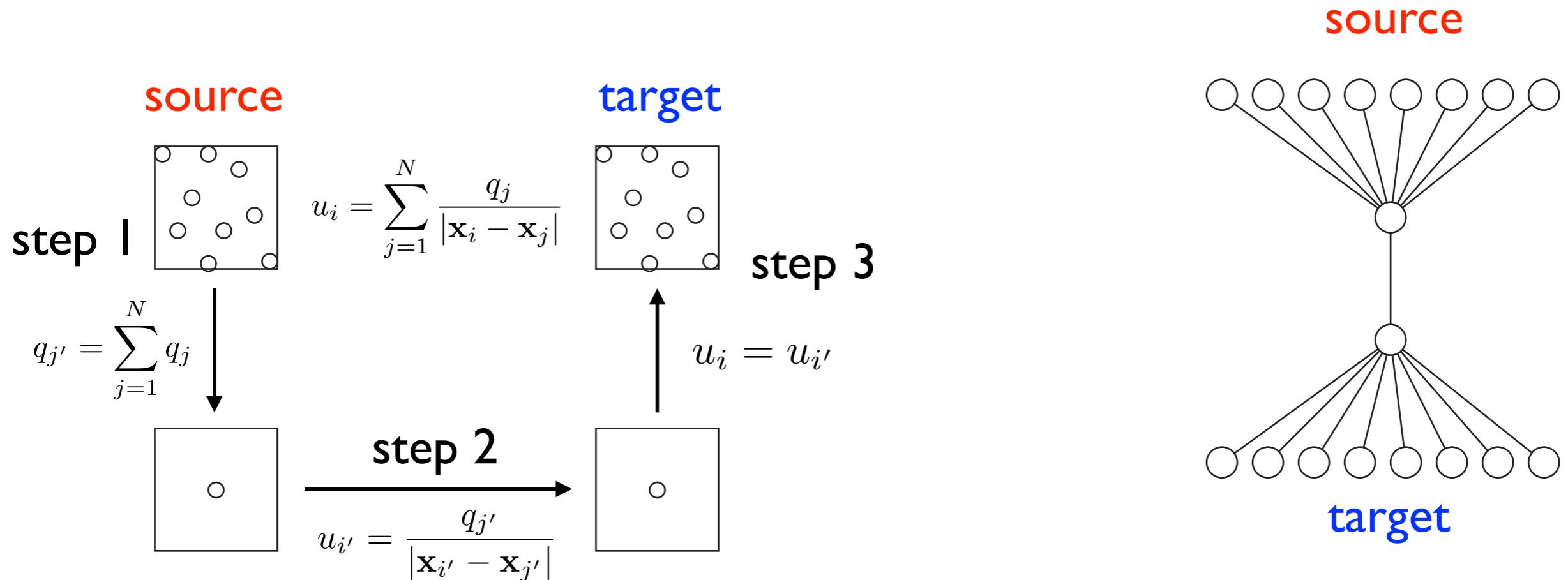
Direct sum - $O(N^2)$

Treecode - $O(N \log N)$

Fast Multipole Method - $O(N)$



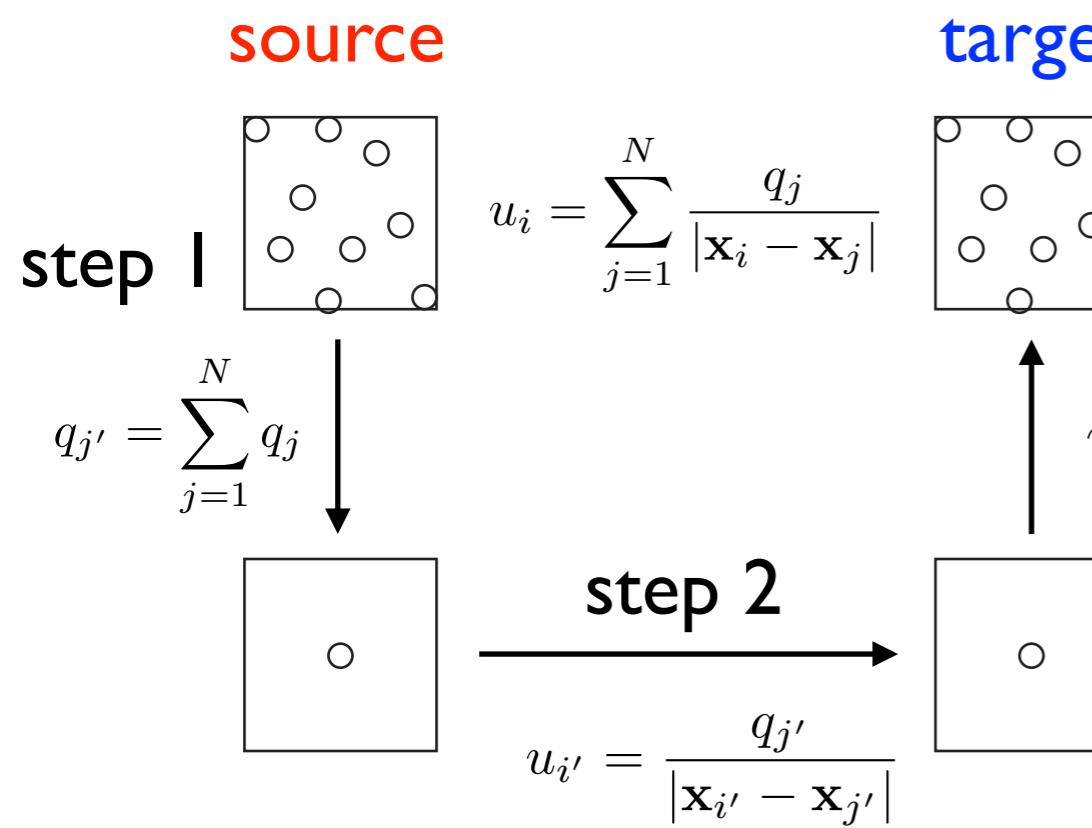
Approximating the interaction



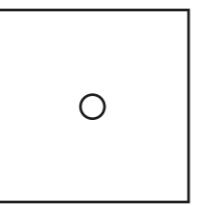
1. Sum all charges
2. Calculate effect of center source on center target
3. Assume all targets in the box have equal potential

Near-far decomposition

non-neighbors



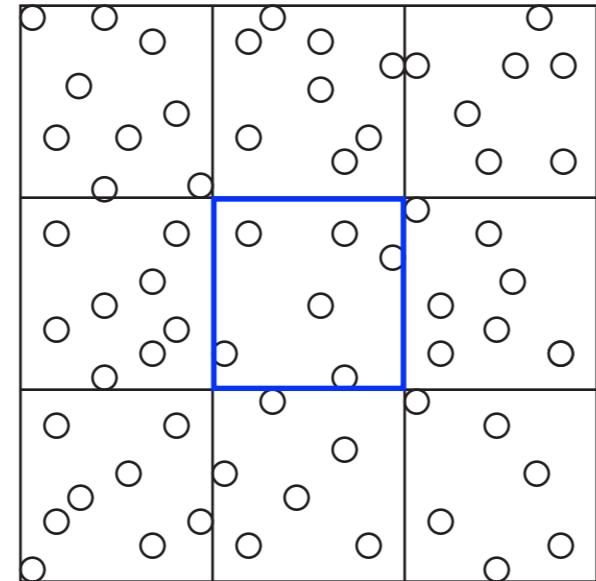
step 3



$$u_i = u_{i'}$$

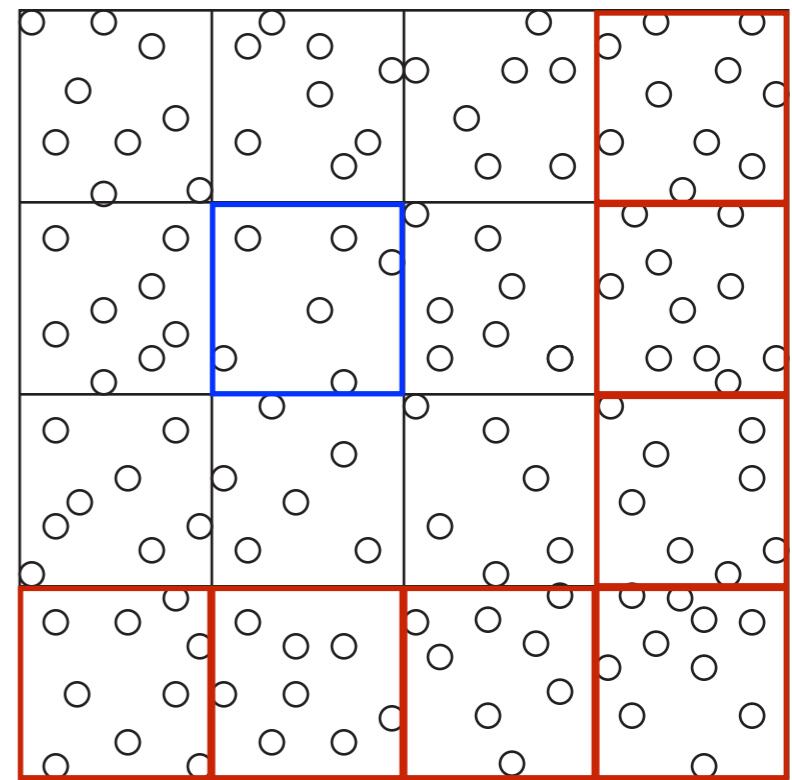
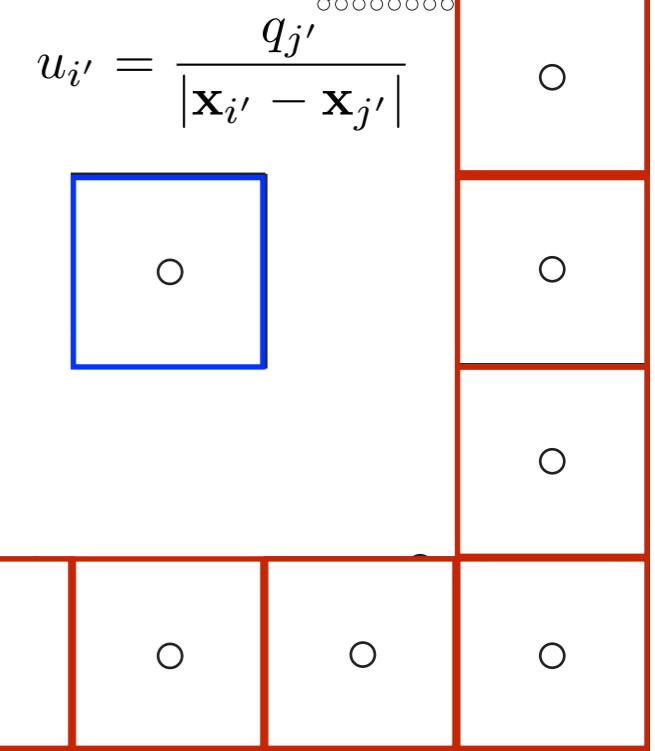
step 2

$$u_{i'} = \frac{q_{j'}}{|\mathbf{x}_{i'} - \mathbf{x}_{j'}|}$$



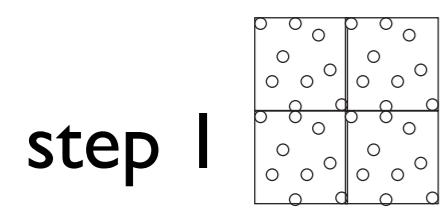
near

$$u_i = \sum_{j=1}^N \frac{q_j}{|\mathbf{x}_i - \mathbf{x}_j|}$$



Hierarchical decomposition

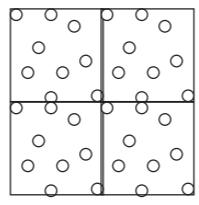
source



step 1

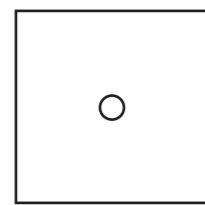
$$q_{j'} = \sum_{j=1}^{N/4} q_j$$

target



step 2

$$q_{j''} = \sum_{j'=1}^4 q_{j'}$$

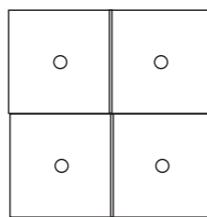


step 3

$$u_{i''} = \frac{q_{j''}}{|\mathbf{x}_{i''} - \mathbf{x}_{j''}|}$$

step 5

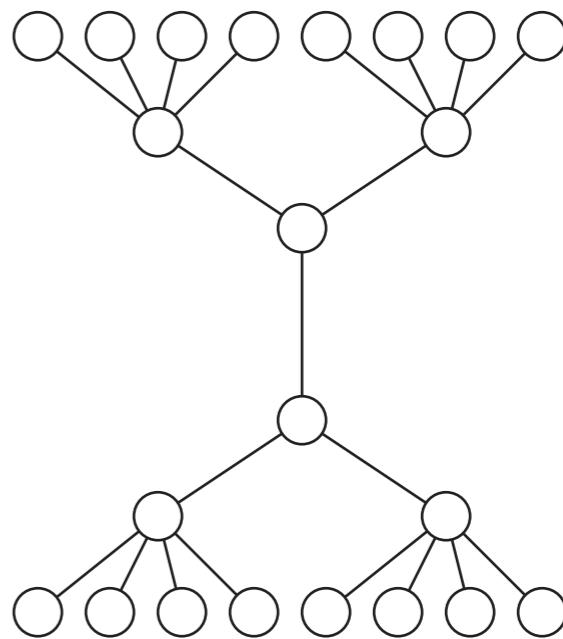
$$u_i = u_{i'}$$



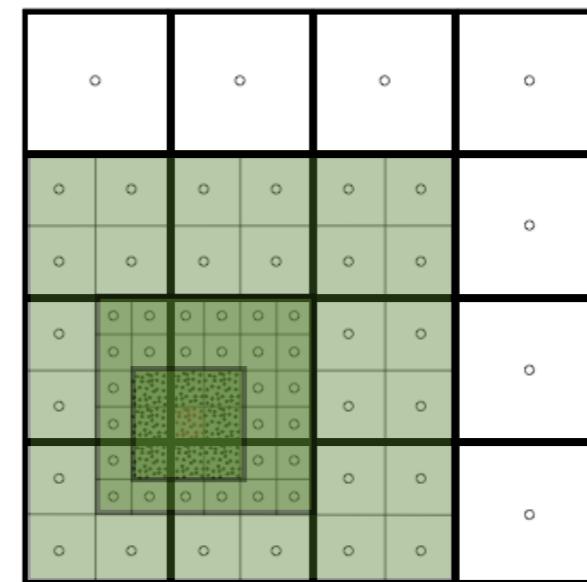
step 4

$$u_{i'} = u_{i''}$$

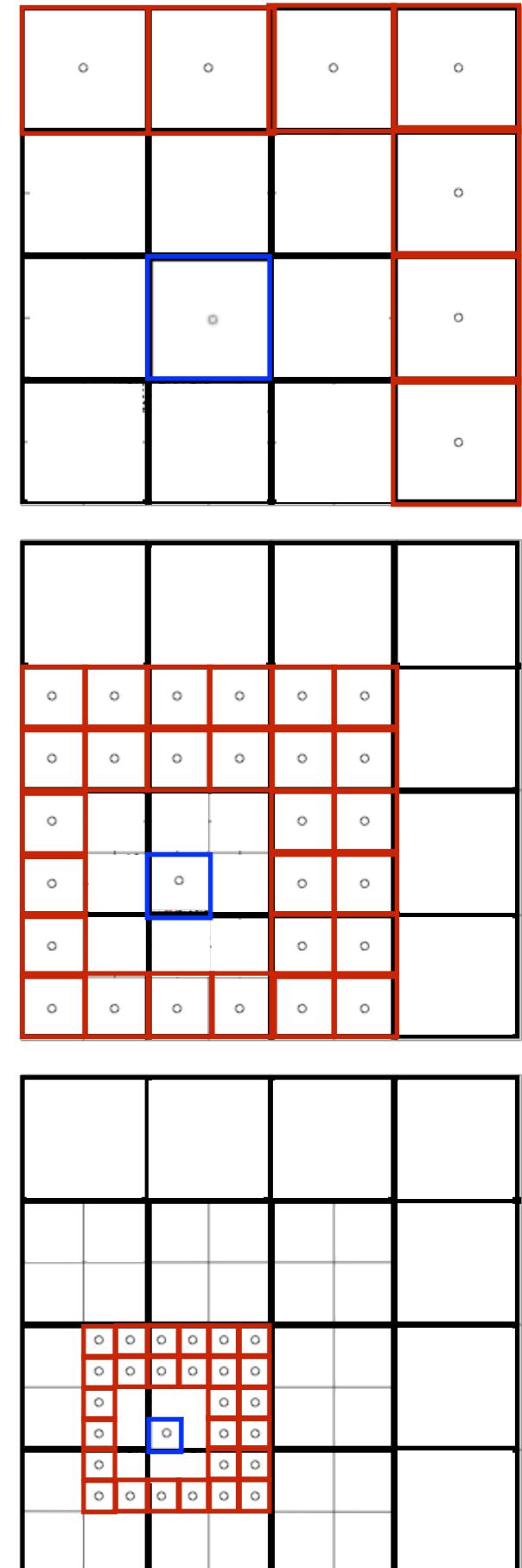
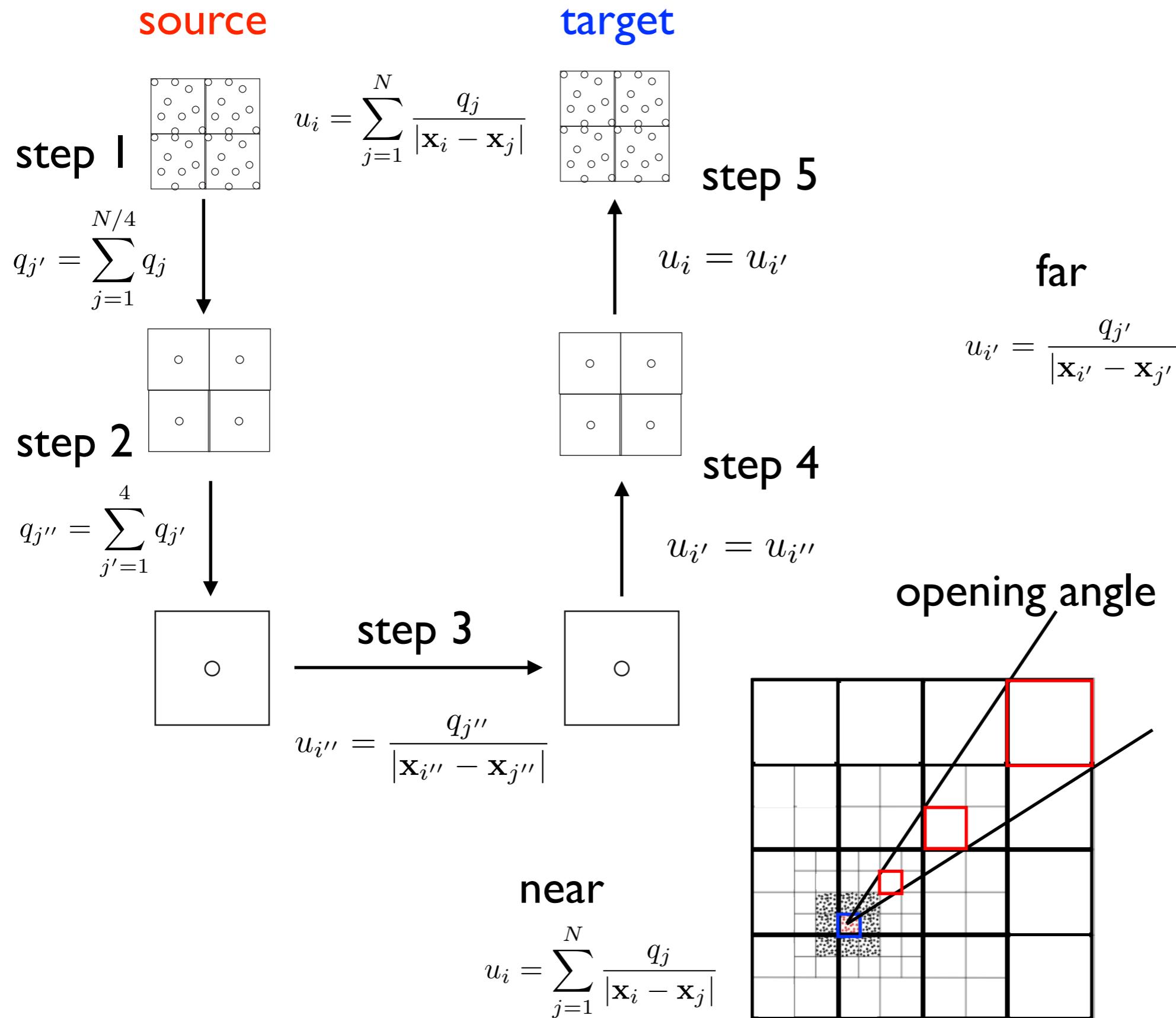
source



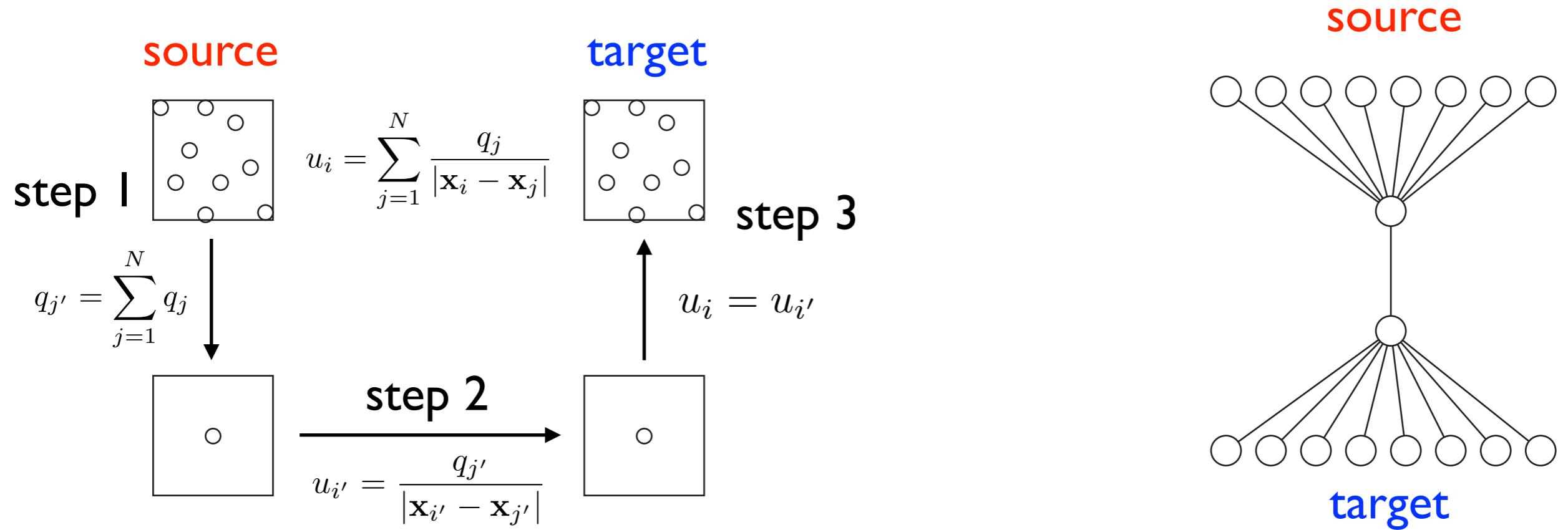
target



Hierarchical near-far decomposition

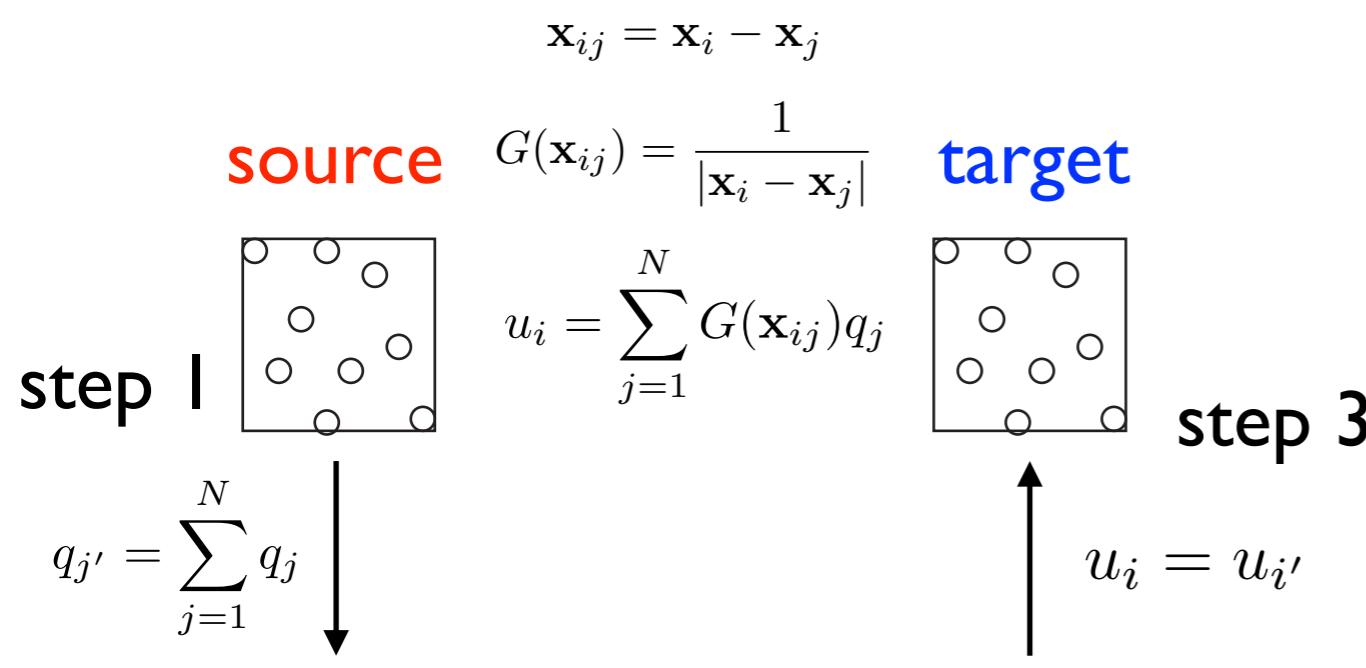


Approximating the interaction



How accurate is the solution?

Higher order approximations



step 2

$$u_{i'} = G(\mathbf{x}_{i'j'}) q_{j'}$$

Binomial theorem

$$(x+y)^n = \sum_{k=0}^n \frac{n!}{(n-k)!k!} x^{n-k} y^k$$

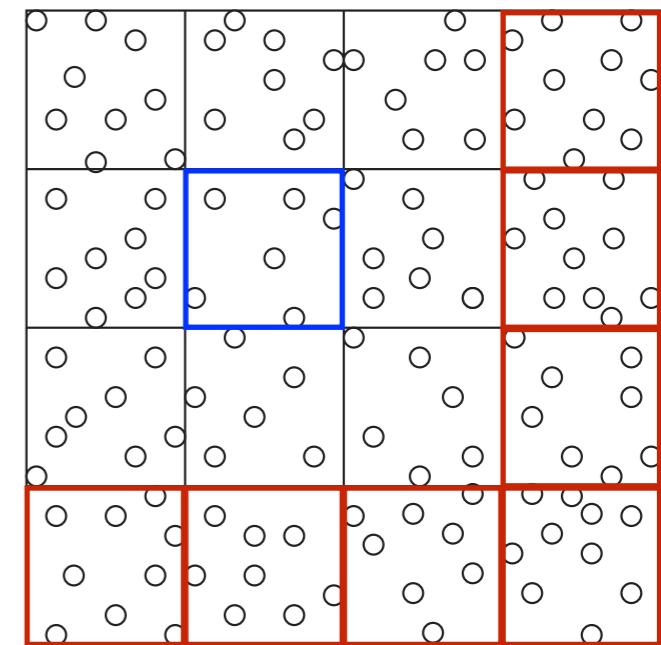


$$\mathbf{x}_{ii'} + \mathbf{x}_{j'j} \ll \mathbf{x}_{i'j'}$$

$$G(\mathbf{x}_{ij}) = \sum_{\mathbf{n}=0}^{\infty} \frac{1}{\mathbf{n}!} (\mathbf{x}_{ii'} + \mathbf{x}_{j'j})^{\mathbf{n}} \nabla^{(\mathbf{n})} G(\mathbf{x}_{i'j'})$$

Taylor expansion

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$



$$G(\mathbf{x}_{ij}) = \sum_{\mathbf{n}=0}^p \frac{1}{\mathbf{n}!} (\mathbf{x}_{ii'} + \mathbf{x}_{j'j})^{\mathbf{n}} \nabla^{(\mathbf{n})} G(\mathbf{x}_{i'j'})$$

$$\rightarrow = \sum_{\mathbf{n}=0}^p \frac{1}{\mathbf{n}!} \sum_{\mathbf{k}=0}^{\mathbf{n}} \frac{\mathbf{n}!}{(\mathbf{n}-\mathbf{k})!\mathbf{k}!} \mathbf{x}_{ii'}^{\mathbf{k}} \mathbf{x}_{j'j}^{\mathbf{n}-\mathbf{k}} \nabla^{(\mathbf{n})} G(\mathbf{x}_{i'j'})$$

Cancel $\mathbf{n}!$ →

$$= \sum_{\mathbf{n}=0}^p \sum_{\mathbf{k}=0}^{\mathbf{n}} \frac{1}{(\mathbf{n}-\mathbf{k})!\mathbf{k}!} \mathbf{x}_{ii'}^{\mathbf{k}} \mathbf{x}_{j'j}^{\mathbf{n}-\mathbf{k}} \nabla^{(\mathbf{n})} G(\mathbf{x}_{i'j'})$$

Swap loop order between \mathbf{n} and \mathbf{k} →

$$= \sum_{\mathbf{k}=0}^p \sum_{\mathbf{n}=\mathbf{k}}^{\mathbf{p}} \frac{1}{(\mathbf{n}-\mathbf{k})!\mathbf{k}!} \mathbf{x}_{ii'}^{\mathbf{k}} \mathbf{x}_{j'j}^{\mathbf{n}-\mathbf{k}} \nabla^{(\mathbf{n})} G(\mathbf{x}_{i'j'})$$

Redefine $\mathbf{n} - \mathbf{k}$ to \mathbf{n} →

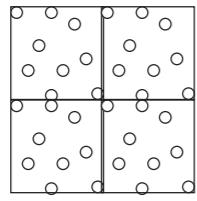
$$= \sum_{\mathbf{k}=0}^p \sum_{\mathbf{n}=0}^{\mathbf{p}-\mathbf{k}} \frac{1}{\mathbf{n}!\mathbf{k}!} \mathbf{x}_{ii'}^{\mathbf{k}} \mathbf{x}_{j'j}^{\mathbf{n}} \nabla^{(\mathbf{n}+\mathbf{k})} G(\mathbf{x}_{i'j'})$$

$$= \sum_{\mathbf{k}=0}^p \frac{1}{\mathbf{k}!} \mathbf{x}_{ii'}^{\mathbf{k}} \underbrace{\sum_{\mathbf{n}=0}^{\mathbf{p}-\mathbf{k}} \nabla^{(\mathbf{n}+\mathbf{k})} G(\mathbf{x}_{i'j'})}_{\mathbf{L}} \underbrace{\frac{1}{\mathbf{n}!} \mathbf{x}_{j'j}^{\mathbf{n}}}_{\mathbf{M}}$$

Multi-level case

source

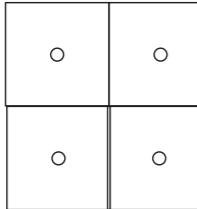
P2M



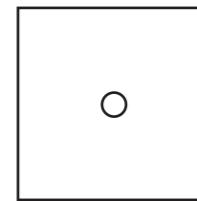
$$u_i = \sum_{j=1}^N G(\mathbf{x}_{ij}) q_j$$

$$\mathbf{M}^{\mathbf{n}}(\mathbf{x}_{j'}) = \sum_{j=1}^{N/4} \frac{1}{\mathbf{n}!} \mathbf{x}_{j'j}^{\mathbf{n}} q_j$$

M2M



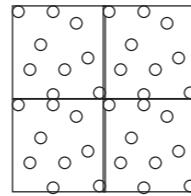
$$\mathbf{M}^{\mathbf{n}}(\mathbf{x}_{j''}) = \sum_{j'=1}^4 \sum_{\mathbf{k}=0}^{\mathbf{n}} \frac{1}{(\mathbf{n}-\mathbf{k})!} \mathbf{x}_{j''j'}^{\mathbf{n}-\mathbf{k}} \mathbf{M}^{\mathbf{k}}(\mathbf{x}_{j'})$$



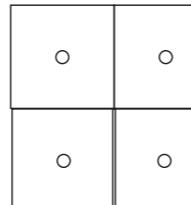
M2L

target

L2P



$$u_i = \sum_{\mathbf{k}=0}^p \frac{1}{\mathbf{k}!} \mathbf{x}_{ii'}^{\mathbf{k}} \mathbf{L}^{\mathbf{k}}(\mathbf{x}_{i'})$$

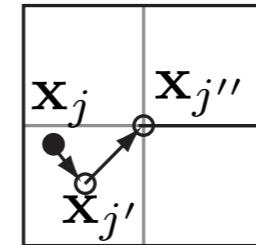


L2L

$$\mathbf{L}^{\mathbf{n}}(\mathbf{x}_{i'}) = \sum_{\mathbf{k}=\mathbf{n}}^p \frac{1}{(\mathbf{k}-\mathbf{n})!} \mathbf{x}_{i'i''}^{\mathbf{k}-\mathbf{n}} \mathbf{L}^{\mathbf{k}}(\mathbf{x}_{i''})$$

$$\mathbf{L}^{\mathbf{k}}(\mathbf{x}_{i''}) = \sum_{\mathbf{n}=0}^{p-\mathbf{k}} \nabla^{(\mathbf{n}+\mathbf{k})} G(\mathbf{x}_{i''j''}) \mathbf{M}^{\mathbf{n}}(\mathbf{x}_{j''})$$

$$\mathbf{L}^{\mathbf{n}}(\mathbf{x}_{i'}) = \sum_{\mathbf{k}=\mathbf{n}}^p \frac{1}{(\mathbf{k}-\mathbf{n})!} \mathbf{x}_{i'i''}^{\mathbf{k}-\mathbf{n}} \mathbf{L}^{\mathbf{k}}(\mathbf{x}_{i''}) \rightarrow = \sum_{\mathbf{n}=0}^p \frac{1}{\mathbf{n}!} \mathbf{x}_{ii'}^{\mathbf{n}} \mathbf{L}^{\mathbf{n}}(\mathbf{x}_{i'})$$



$$G(\mathbf{x}_{ij}) = \underbrace{\sum_{\mathbf{k}=0}^p \frac{1}{\mathbf{k}!} \mathbf{x}_{ii''}^{\mathbf{k}} \sum_{\mathbf{n}=0}^{p-\mathbf{k}} \nabla^{(\mathbf{n}+\mathbf{k})} G(\mathbf{x}_{i''j''})}_{\mathbf{L}} \underbrace{\frac{1}{\mathbf{n}!} \mathbf{x}_{j''j}^{\mathbf{n}}}_{\mathbf{M}}$$

$$\mathbf{x}_{j''j} = \mathbf{x}_{j''j'} + \mathbf{x}_{j'j}$$

$$\mathbf{M}^{\mathbf{n}}(\mathbf{x}_{j''}) = \frac{1}{\mathbf{n}!} (\mathbf{x}_{j''j'} + \mathbf{x}_{j'j})^{\mathbf{n}} q_j$$

$$= \frac{1}{\mathbf{n}!} \sum_{\mathbf{k}=0}^{\mathbf{n}} \frac{\mathbf{n}!}{(\mathbf{n}-\mathbf{k})!\mathbf{k}!} \mathbf{x}_{j''j'}^{\mathbf{n}-\mathbf{k}} \mathbf{x}_{j'j}^{\mathbf{k}} q_j$$

$$= \sum_{\mathbf{k}=0}^{\mathbf{n}} \frac{1}{(\mathbf{n}-\mathbf{k})!} \mathbf{x}_{j''j'}^{\mathbf{n}-\mathbf{k}} \mathbf{M}^{\mathbf{k}}(\mathbf{x}_{j'})$$

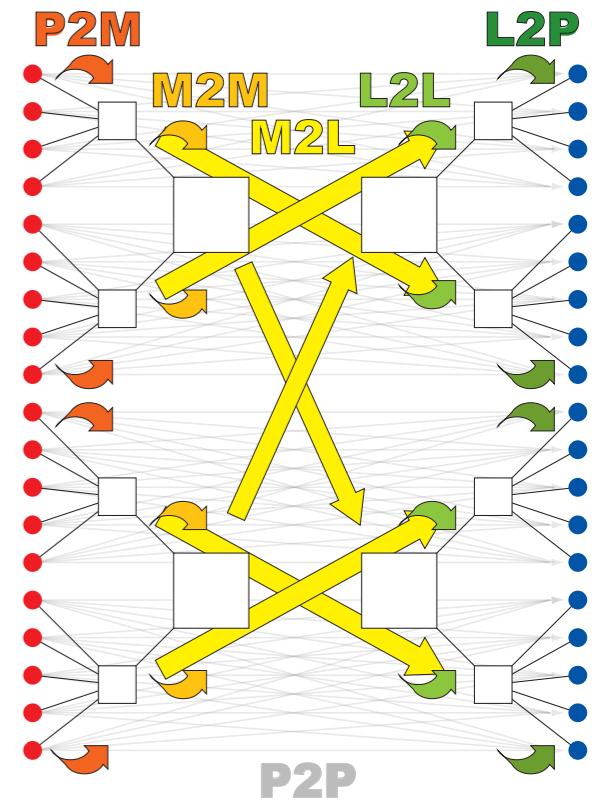
$$\mathbf{x}_{ii''} = \mathbf{x}_{ii'} + \mathbf{x}_{i'i''}$$

$$u_i = \sum_{\mathbf{k}=0}^p \frac{1}{\mathbf{k}!} (\mathbf{x}_{ii'} + \mathbf{x}_{i'i''})^{\mathbf{k}} \mathbf{L}^{\mathbf{k}}(\mathbf{x}_{i''})$$

$$= \sum_{\mathbf{k}=0}^p \sum_{\mathbf{n}=0}^{\mathbf{k}} \frac{1}{(\mathbf{k}-\mathbf{n})!\mathbf{n}!} \mathbf{x}_{ii'}^{\mathbf{n}} \mathbf{x}_{i'i''}^{\mathbf{k}-\mathbf{n}} \mathbf{L}^{\mathbf{k}}(\mathbf{x}_{i''})$$

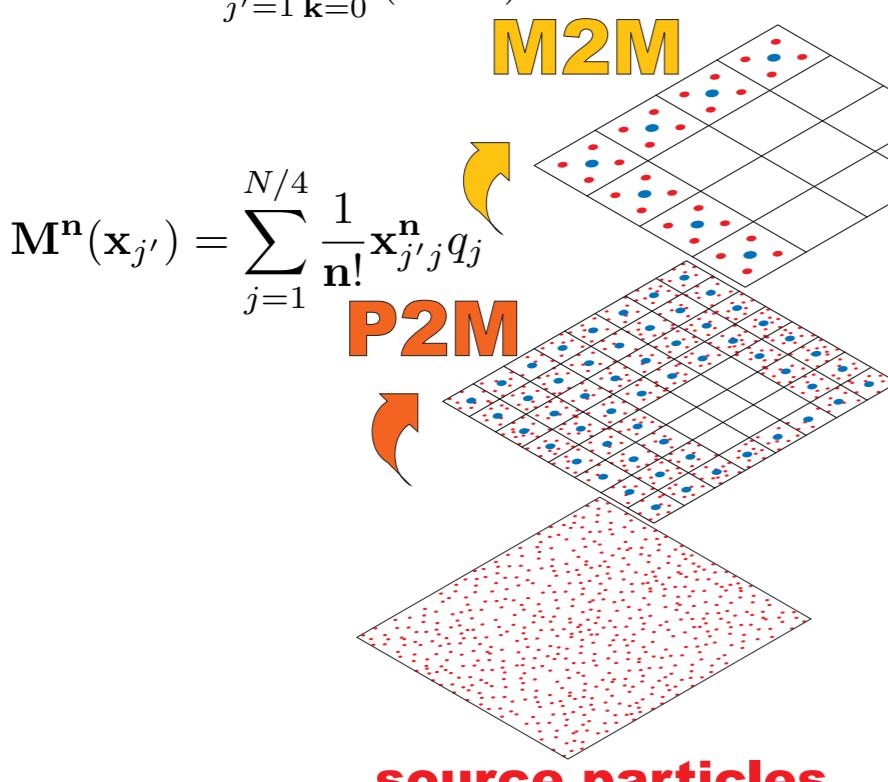
$$= \sum_{\mathbf{n}=0}^p \sum_{\mathbf{k}=\mathbf{n}}^p \frac{1}{(\mathbf{k}-\mathbf{n})!\mathbf{n}!} \mathbf{x}_{ii'}^{\mathbf{n}} \mathbf{x}_{i'i''}^{\mathbf{k}-\mathbf{n}} \mathbf{L}^{\mathbf{k}}(\mathbf{x}_{i''})$$

Flow of Calculation

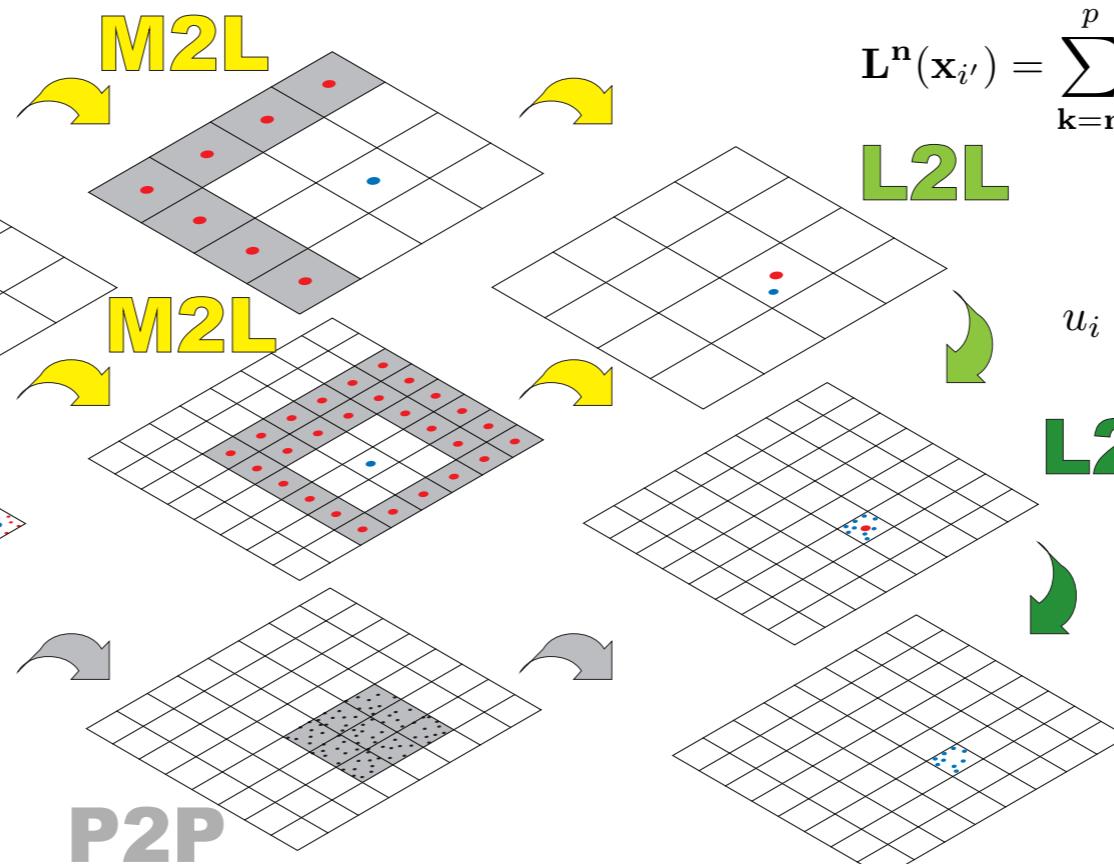


$$L^k(x_{i''}) = \sum_{n=0}^{p-k} \nabla^{(n+k)} G(x_{i''j''}) M^n(x_{j''})$$

$$M^n(x_{j''}) = \sum_{j'=1}^4 \sum_{k=0}^n \frac{1}{(n-k)!} x_{j''j'}^{n-k} M^k(x_{j'})$$

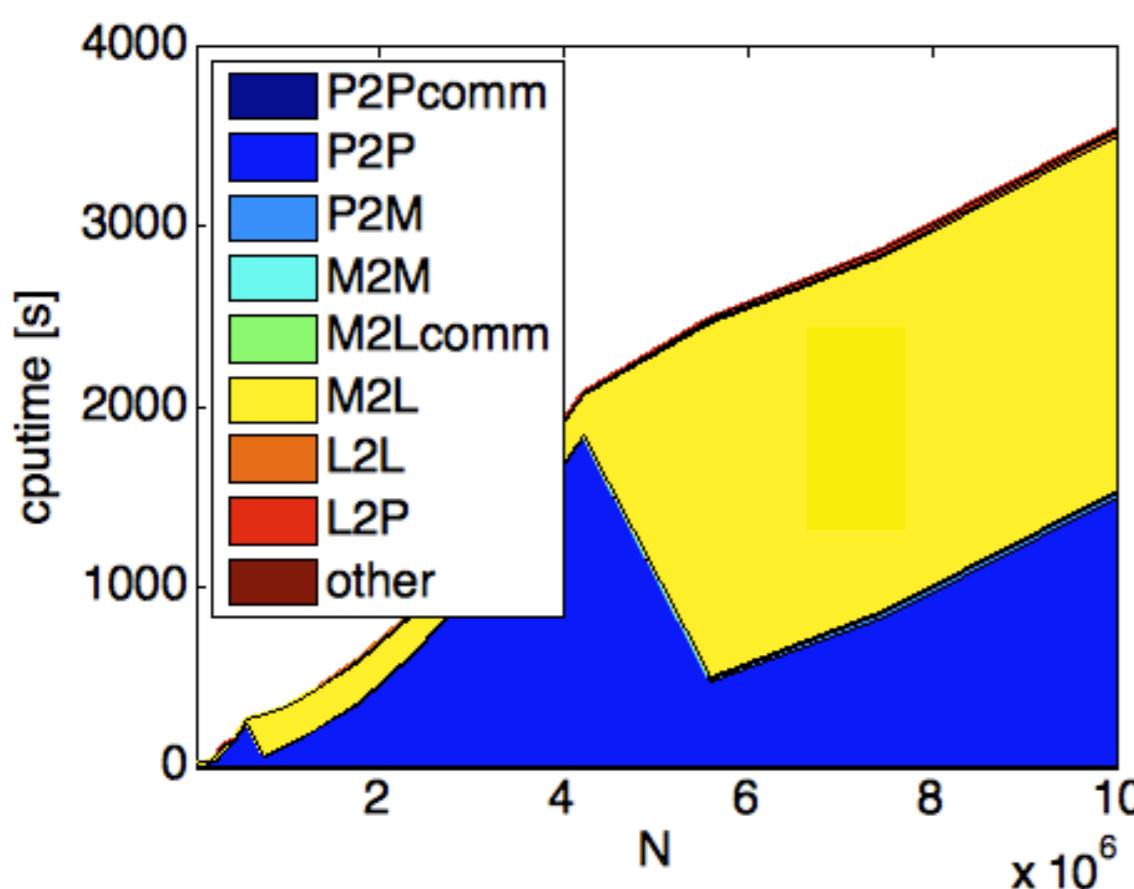
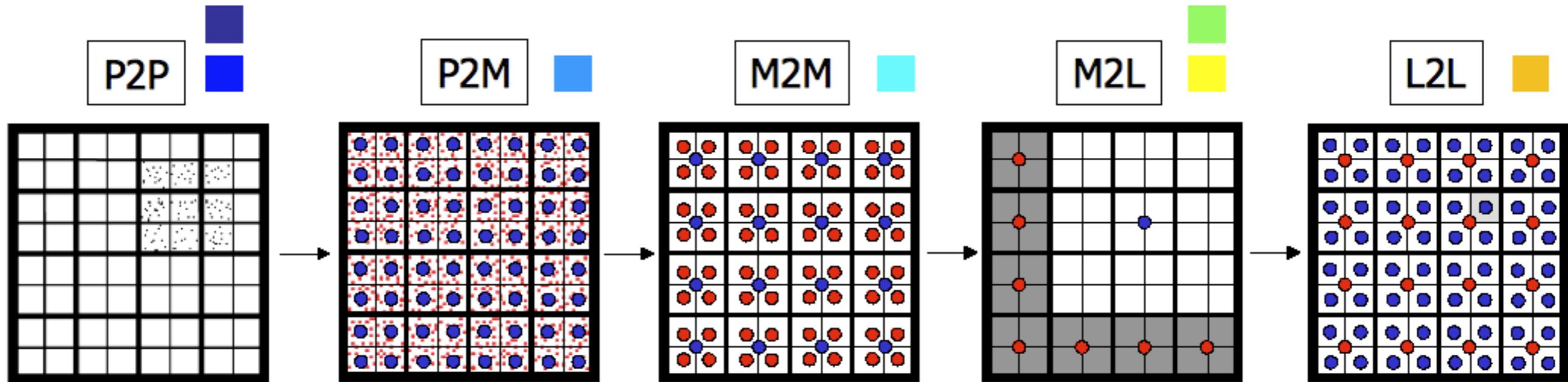


$$M^n(x_{j'}) = \sum_{j=1}^{N/4} \frac{1}{n!} x_{j'j}^n q_j$$



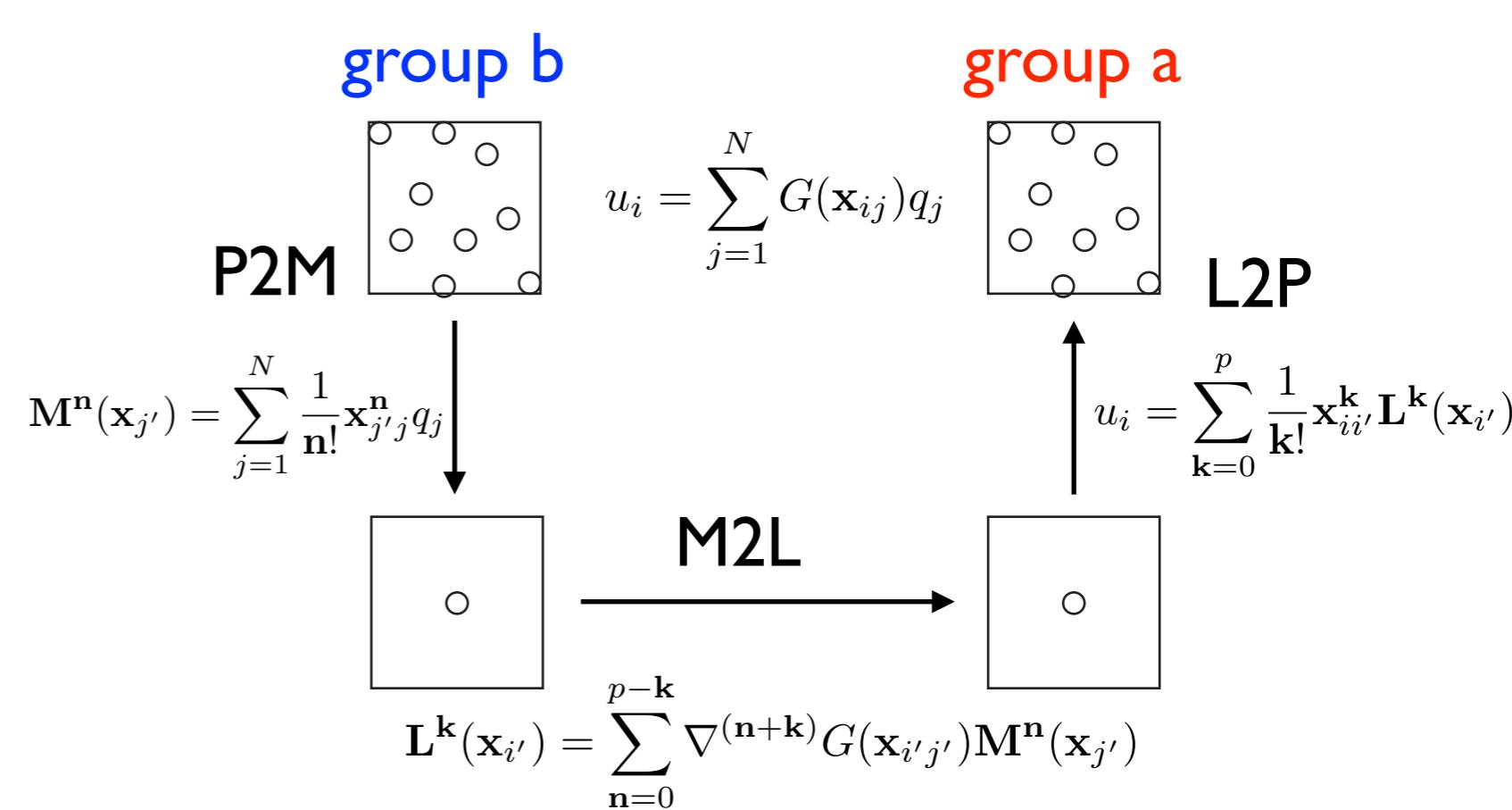
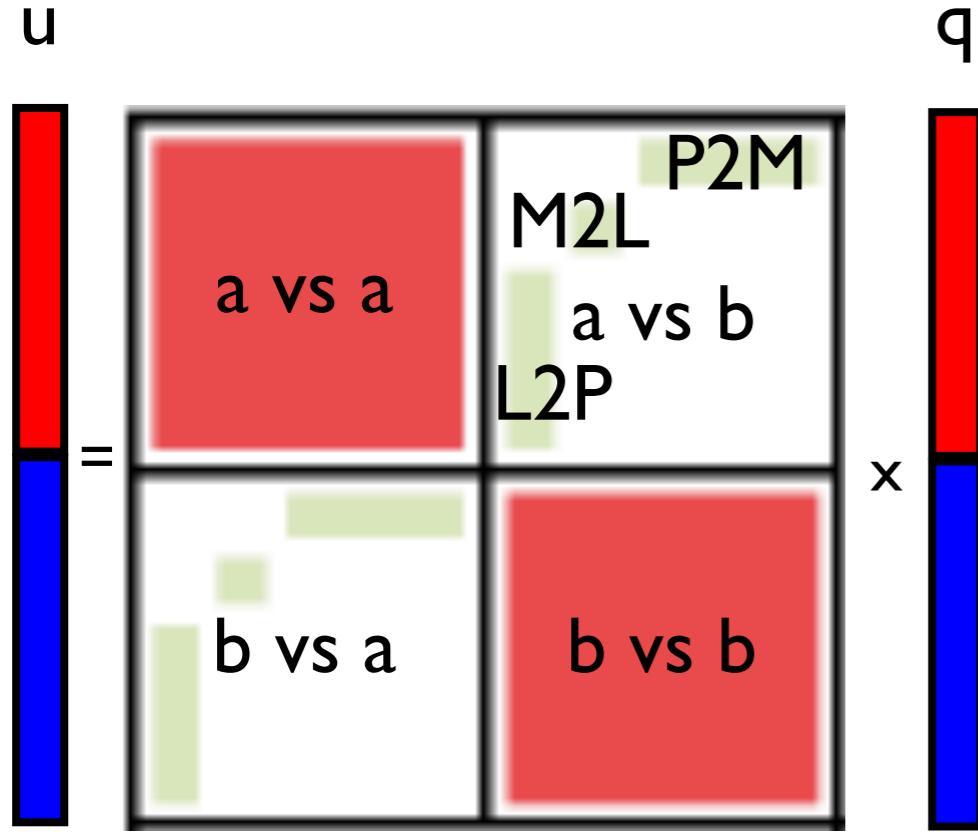
$$u_i = \sum_{j=1}^N G(x_{ij}) q_j$$

How much time does each part take?



Hierarchical low-rank matrices

FMM as a hierarchical matrix-vector multiplication



Consider the case $p=3$

$$\begin{bmatrix} M_0 \\ M_1 \\ M_2 \end{bmatrix} = \begin{bmatrix} P2M & \frac{1}{n!} \mathbf{x}_{j'j}^n \end{bmatrix} \times \begin{bmatrix} q_0 \\ q_1 \\ q_2 \\ q_3 \\ q_4 \\ q_5 \\ q_6 \\ q_7 \\ q_8 \end{bmatrix}$$

$$\begin{bmatrix} L_0 \\ L_1 \\ L_2 \end{bmatrix} = \begin{bmatrix} M2L \\ G(\mathbf{x}_{i'j'}) \end{bmatrix} \times \begin{bmatrix} M_0 \\ M_1 \\ M_2 \end{bmatrix}$$

$$\begin{bmatrix} u_0 \\ u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \\ u_6 \\ u_7 \\ u_8 \end{bmatrix} = \begin{bmatrix} L2P \\ \frac{1}{k!} \mathbf{x}_{ii'}^k \end{bmatrix}$$

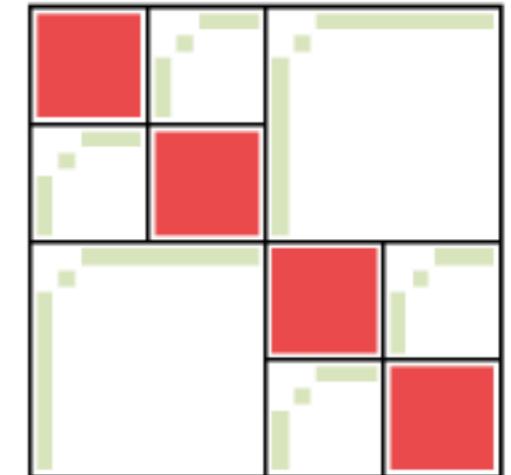
Hierarchical low-rank matrices

Replace dense linear algebra

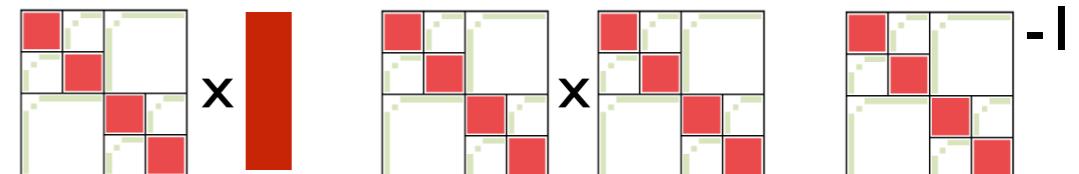
Compute : $\mathcal{O}(N^3) \longrightarrow \mathcal{O}(N)$

Memory : $\mathcal{O}(N^2) \longrightarrow \mathcal{O}(N)$

Hierarchical off-diagonal blocks
Approximated with low rank



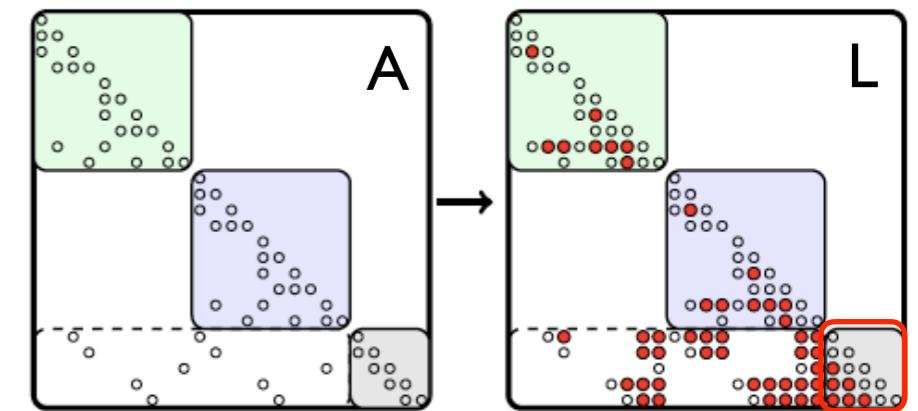
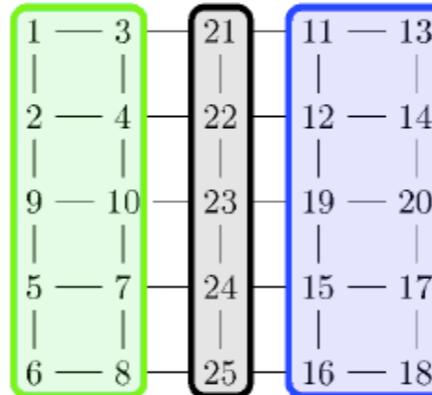
Augment sparse linear algebra



Sparse direct solvers

Schur complement (frontal matrix) is dense but numerically low-rank

Nested dissection



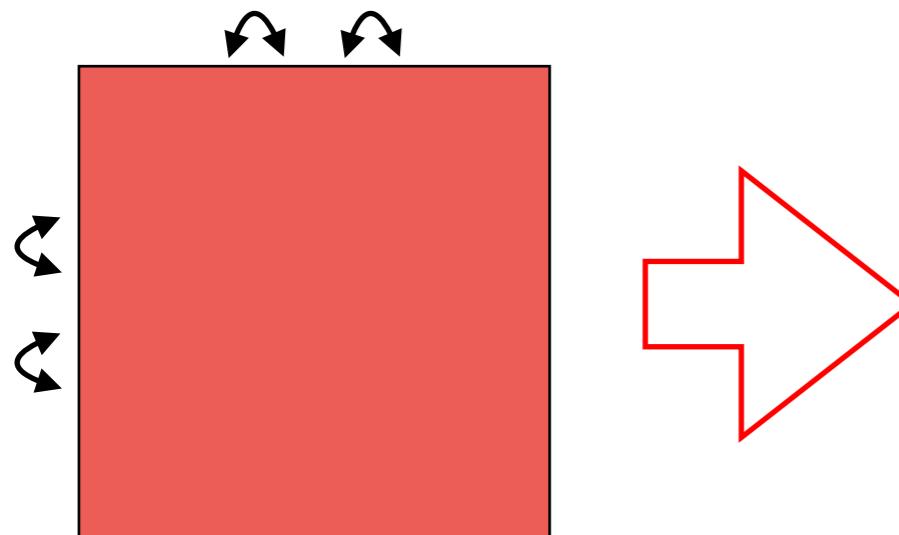
Iterative solvers

Use small rank to precondition

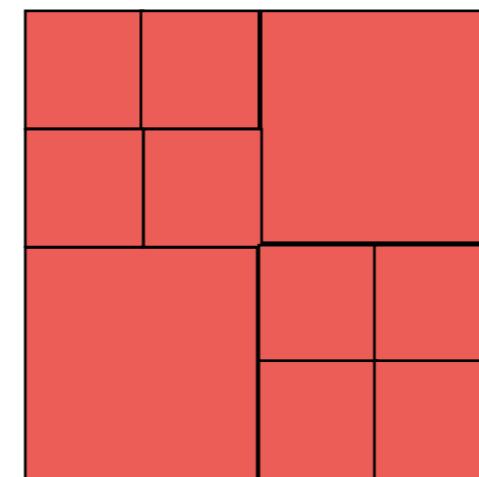
Less sensitive to matrix condition than multigrid

Three Stages of H-matrix Compression

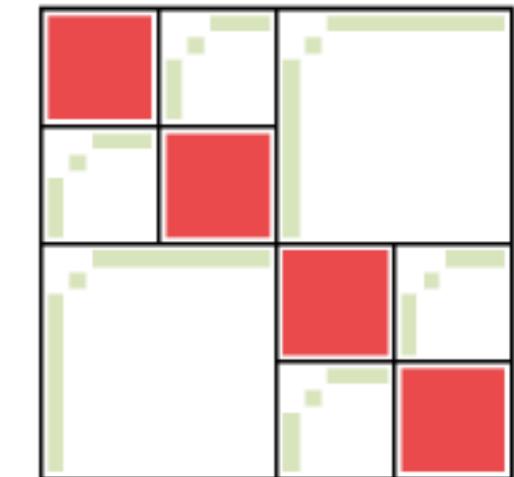
Reordering



Subdivision



Low-Rank



What to minimize ?

Fill-in (Graph connections)
Rank (Geometric distance)
Communication (Locality)
→ Close nodes are usually connected, so minimizing rank will minimize fill-in

How much to divide ?

Subdividing the block will decrease the rank

The rank can be kept constant while using the subdivision to control the accuracy
→ SIMD friendly

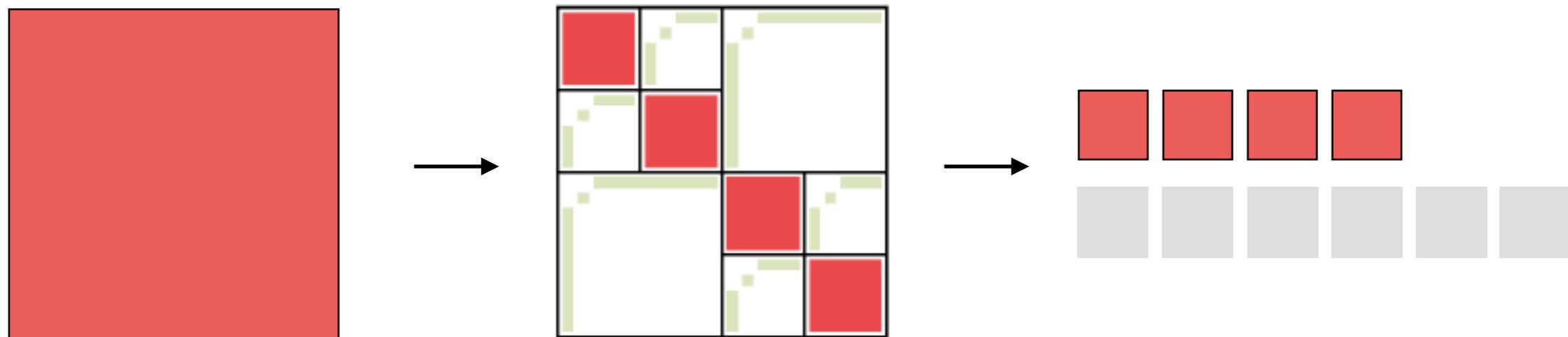
Speed or reliability ?

ACA is fast but unreliable
RSVD is reliable but slow

Many small RSVDs must be accelerated
→ TSQR on Tensor Cores

These methods run efficiently on modern architectures

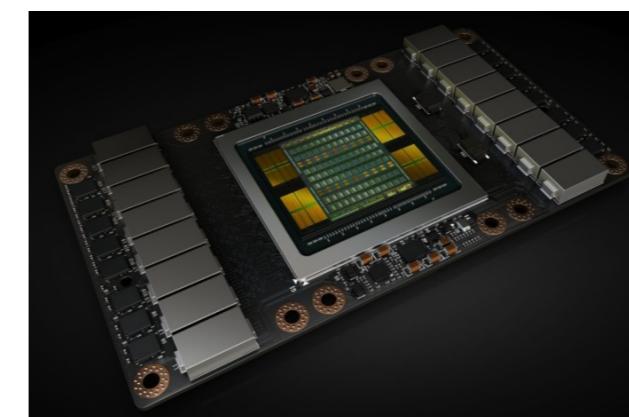
Batch of many small dense matrices



Low-rank approximation needs low arithmetic precision

$$D = \left(\begin{array}{cccc} A_{0,0} & A_{0,1} & A_{0,2} & A_{0,3} \\ A_{1,0} & A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,0} & A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,0} & A_{3,1} & A_{3,2} & A_{3,3} \end{array} \right) \text{FP16 or FP32} \left(\begin{array}{cccc} B_{0,0} & B_{0,1} & B_{0,2} & B_{0,3} \\ B_{1,0} & B_{1,1} & B_{1,2} & B_{1,3} \\ B_{2,0} & B_{2,1} & B_{2,2} & B_{2,3} \\ B_{3,0} & B_{3,1} & B_{3,2} & B_{3,3} \end{array} \right) \text{FP16} + \left(\begin{array}{cccc} C_{0,0} & C_{0,1} & C_{0,2} & C_{0,3} \\ C_{1,0} & C_{1,1} & C_{1,2} & C_{1,3} \\ C_{2,0} & C_{2,1} & C_{2,2} & C_{2,3} \\ C_{3,0} & C_{3,1} & C_{3,2} & C_{3,3} \end{array} \right) \text{FP16 or FP32}$$

4 4



Replacing Exact Linear Algebra with Low-Rank

Exact

$$\mathcal{O}(N^3)$$

Approximate

$$\mathcal{O}(N)$$

Application

ScaLAPACK

cuSolverMG

LAPACK

PLASMA

BLAS

CPU

FP64

cuSolverDN

MAGMA

CUBLAS

MKL

Distributed

QR

LU

MatMul

Mat-vec

App.

HiCMA

STRUMPACK

GOFMM

LoRaSp

HBLAS

?PU

TF32, bfloat16



List of implementations

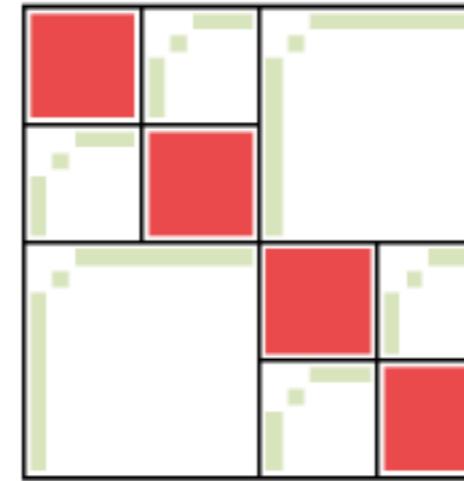
	Method	Developer	url
AHMED	H-matrix	M. Bebendorf	https://github.com/xantares/ahmed
ASKIT	FMM	C. D. Yu	http://padas.ices.utexas.edu/libaskit
DMHM	H-matrix	J. Poulson	https://bitbucket.org/poulson/dmhm/src/default/
GOFMM	H^2 -matrix	C. D. Yu	https://github.com/ChenhanYu/hmlp
H2Lib	H^2 -matrix	S. Börm	https://github.com/H2Lib/H2Lib
H2Tools	H^2 -matrix	A. Mikhalev	https://bitbucket.org/muxas/h2tools
HACApK	H-matrix	A. Ida	https://github.com/HLRA-JHPCN/HACApK-MAGMA
HiCMA	H-matrix	H. Ltaief	https://github.com/ecrc/hicma
HLib	H-matrix	L. Grasedyck	http://www.hlib.org
HLibPro	H-matrix	R. Kriemann	http://www.hlibpro.com
hmglib	H-matrix	P. Zaspel	https://github.com/zaspel/hmglib
HODLR	HODLR	A. Aminfar	https://github.com/amiraa127/Dense_HODLR
HSS	HSS	J. Xia	http://www.math.purdue.edu/~xiaj/
LoRaSp	H^2 -matrix	H. Pouransari	https://bitbucket.org/hadip/lorasp
MUMPS-BLR	BLR	P. R. Amestoy	http://mumps.enseeiht.fr
STURMPACK	HSS	P. Ghysels	http://portal.nersc.gov/project/sparse/strumpack

Differences between LRA methods

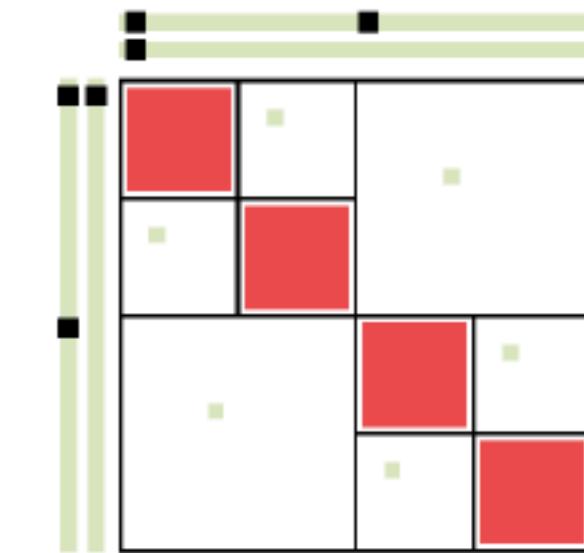
	Shared Basis	Admissibility
H-matrix	No	Strong
H^2 -matrix	No	Strong
HODLR	No	Weak
HSS	Yes	Weak
RS/HIF	Yes	Strong
IFMM	Yes	FMM
(inv)-ASKIT	Yes	Strong
BLR	No	non-hierarchical
BLR ²	Yes	non-hierarchical

Nested Basis

Non-nested

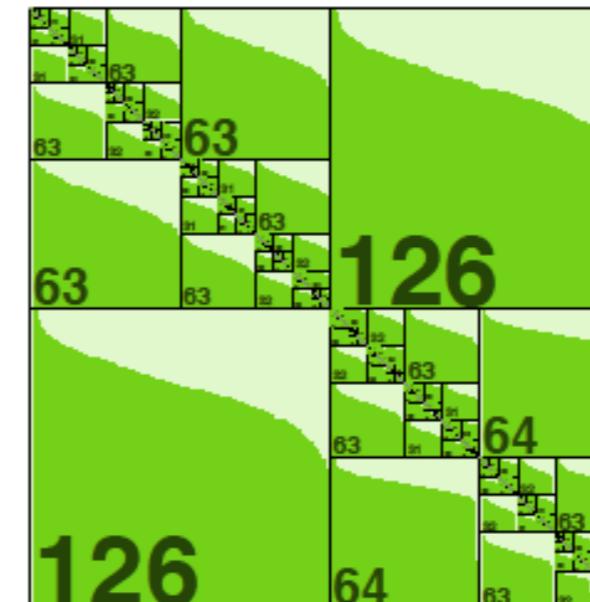


Nested

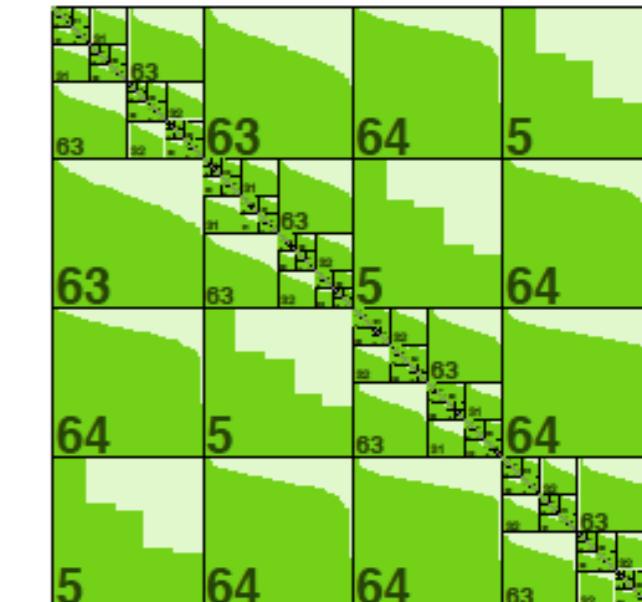


Admissibility

Weak

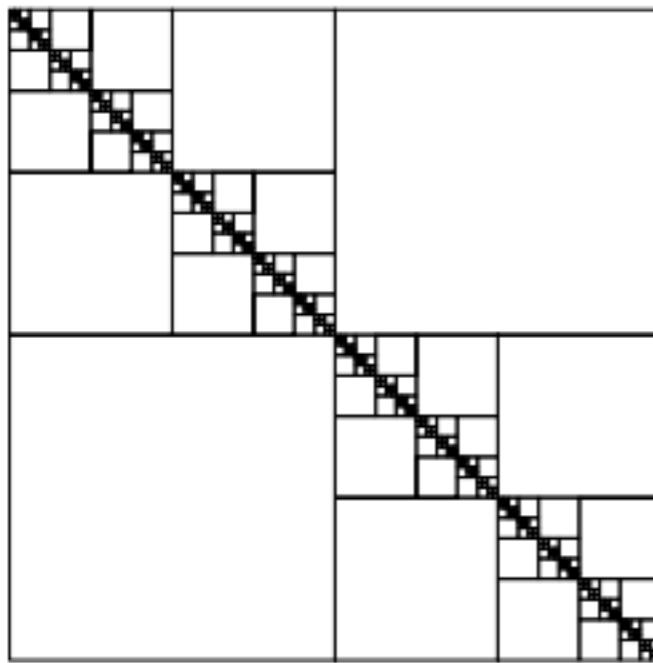


well separated Strong

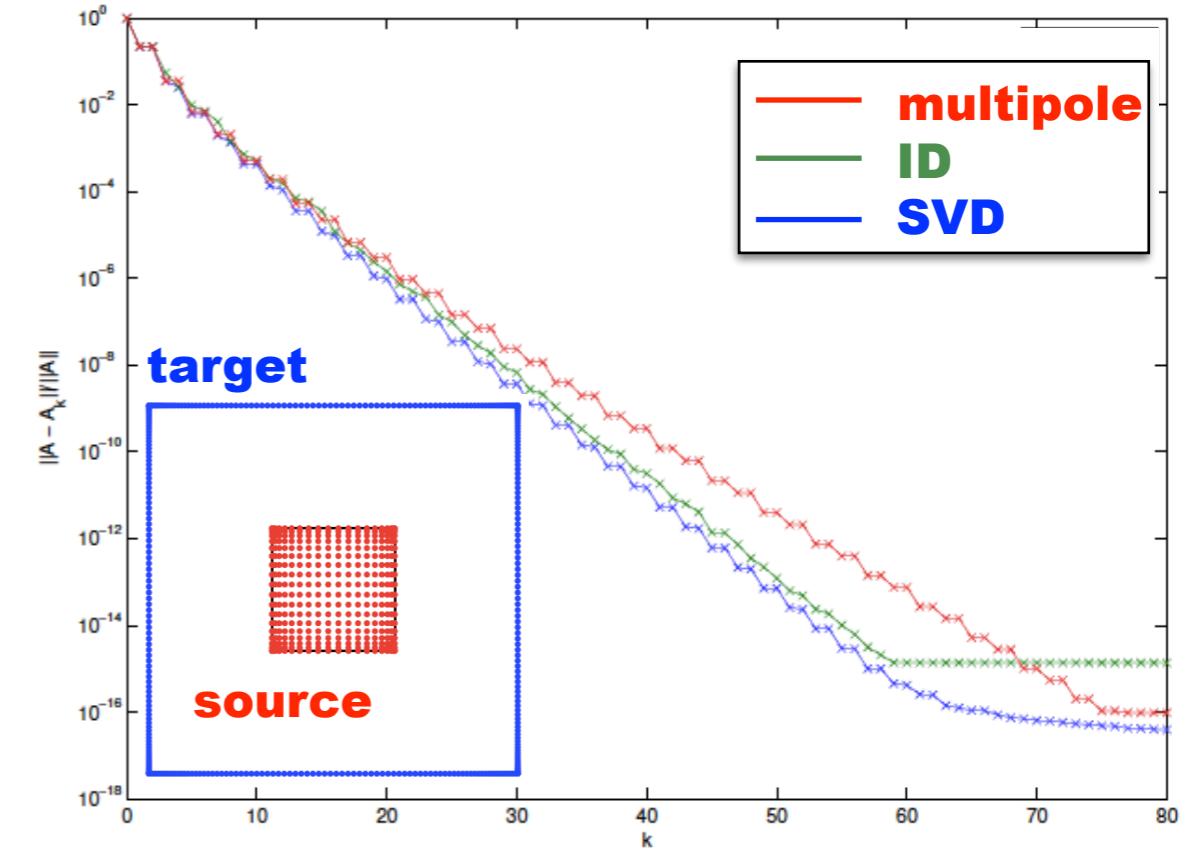
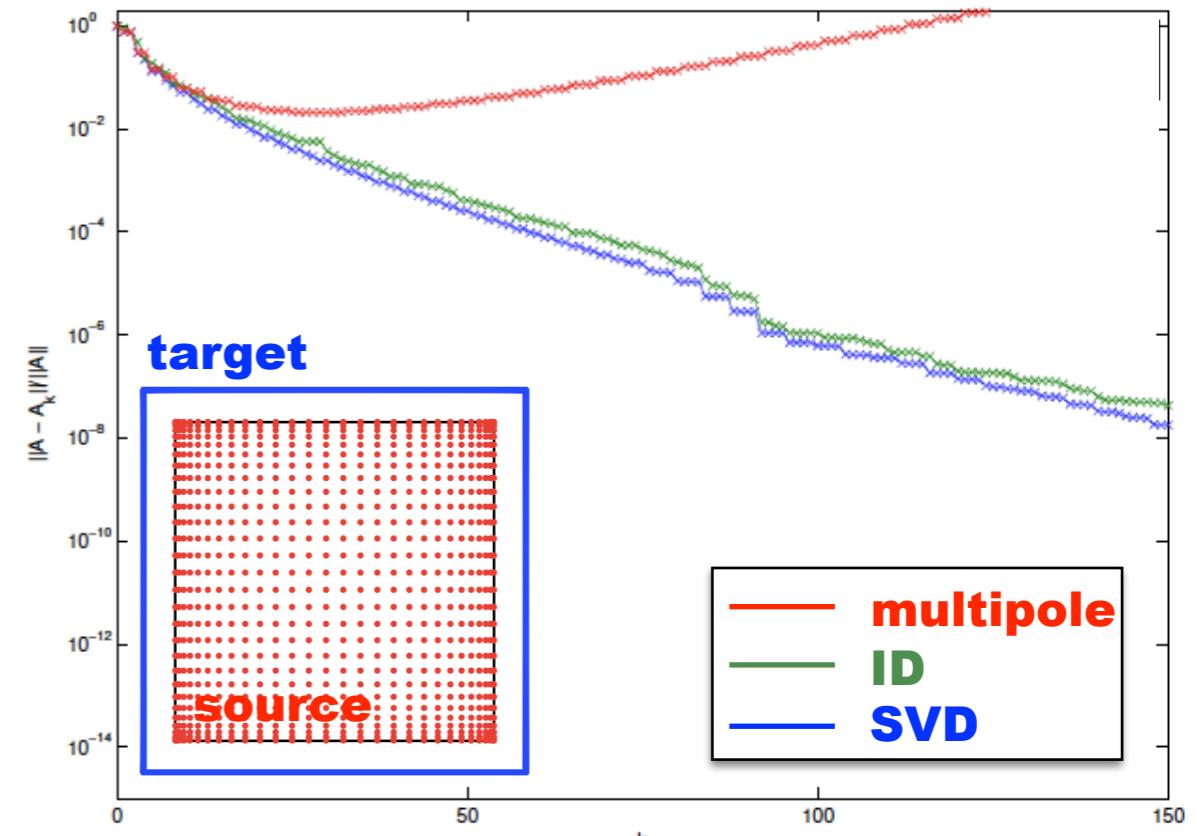
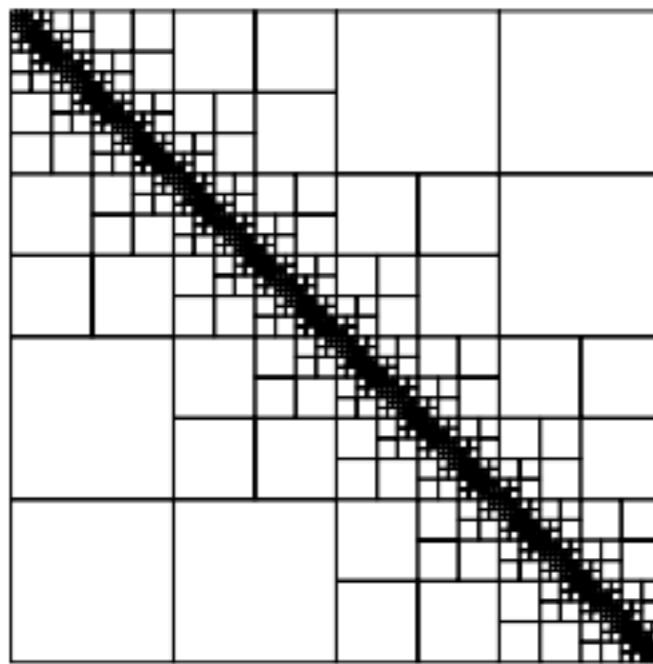


Admissibility condition

Weak admissibility



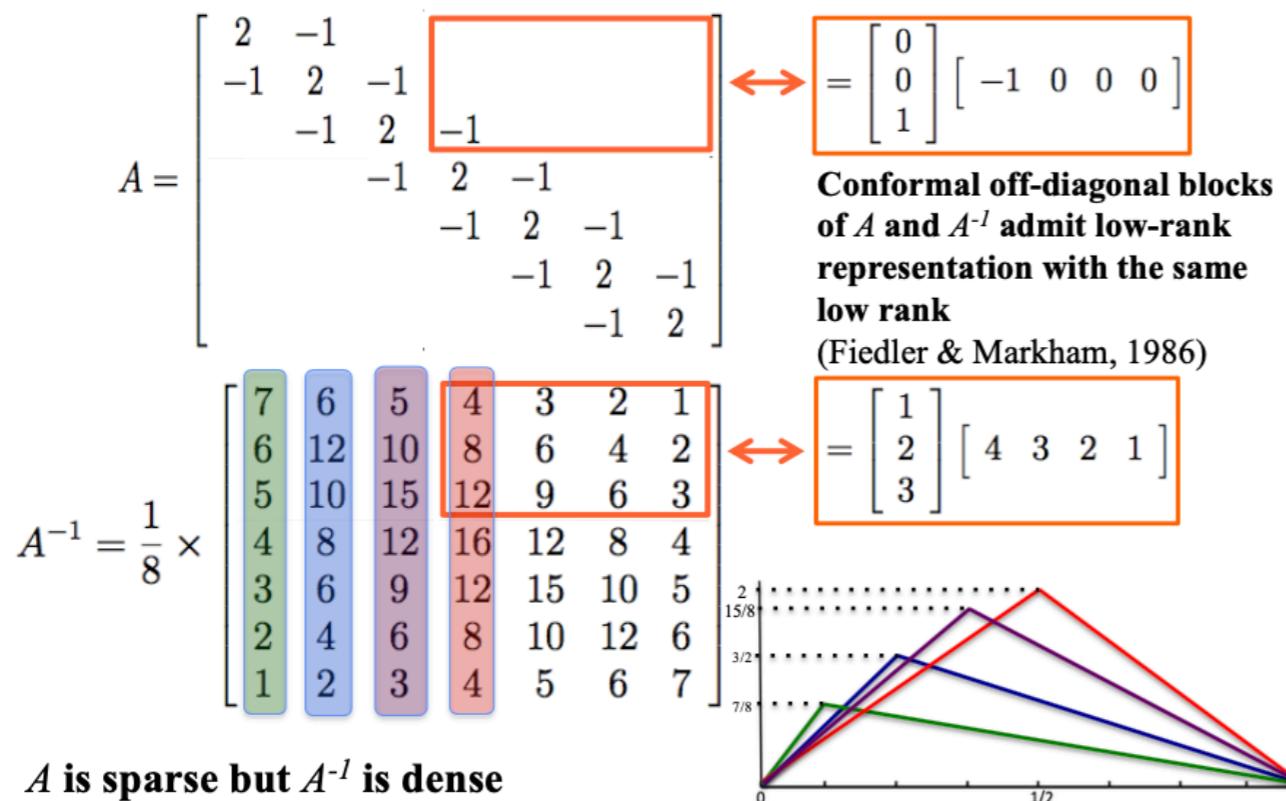
Standard admissibility



Nullity Theorem

David Keyes' lecture 2

Simple example of data sparsity with the 1D Laplacian



$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} E & F \\ G & H \end{bmatrix}$$

nullity $A =$ nullity H ,
nullity $B =$ nullity F ,
nullity $C =$ nullity G ,
nullity $D =$ nullity E .

$$\text{rank}(A) + \text{nullity}(A) = n.$$



↓ Apply it recursively



Nullity Theorem ?

Decay of singular values

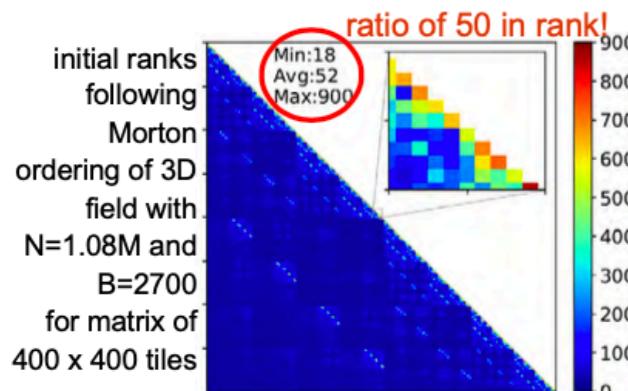
David Keyes' lecture 2

Rank distribution challenges with 3D exponential kernels

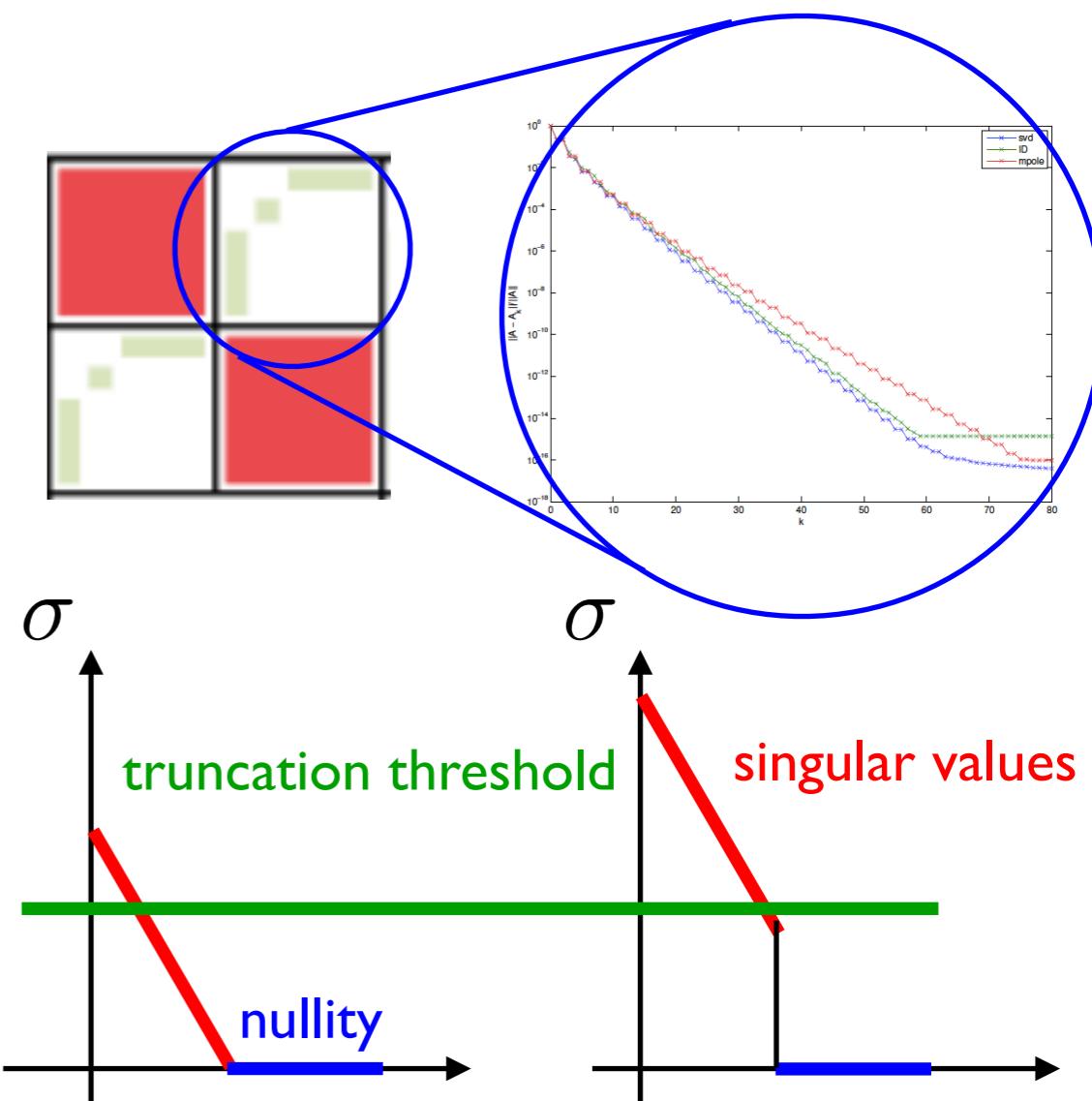
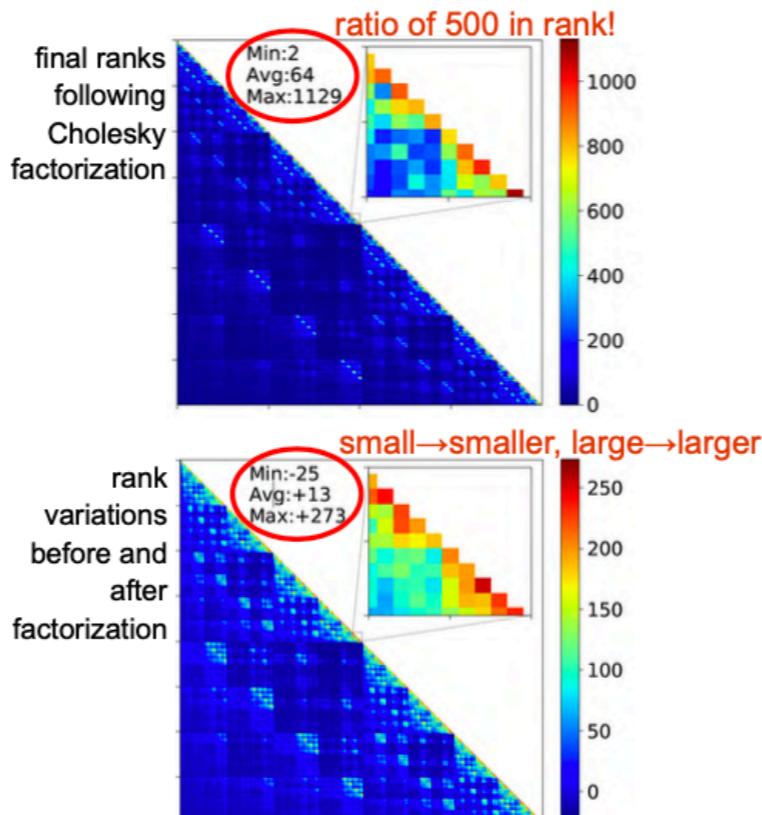
The simple exponential kernel:

$$C(r; \ell) = \exp\left(-\frac{r}{\ell}\right)$$

is suited for rough correlations such as the variation of wind speed or temperature with altitude, and leads to wide rank disparities



Cao, Pei, Akbudak, Bosilca, Ltaief, K. & Dongarra, Leveraging ParSEC Runtime Support to Tackle Challenging 3D Data-sparse Matrix Problems. IPDPS (IEEE), 2021



A

$A * 100$

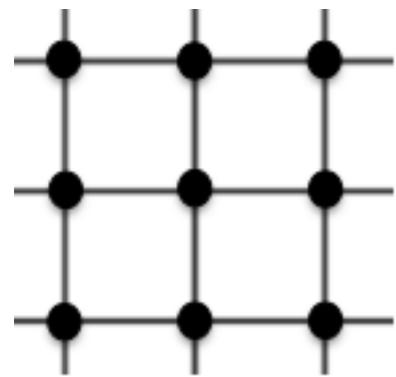
The exact rank does not grow, but the numerical/truncated rank does

Kronecker factorization

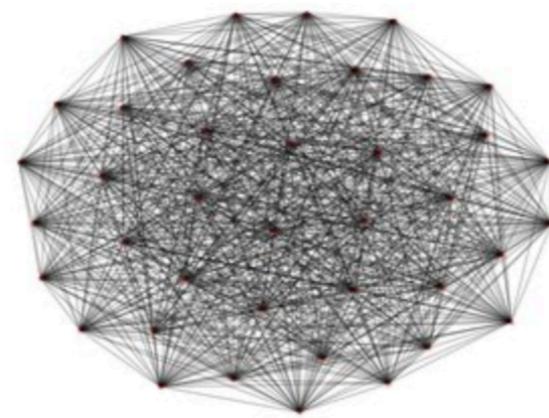
Structure of matrices

2-D or 3-D

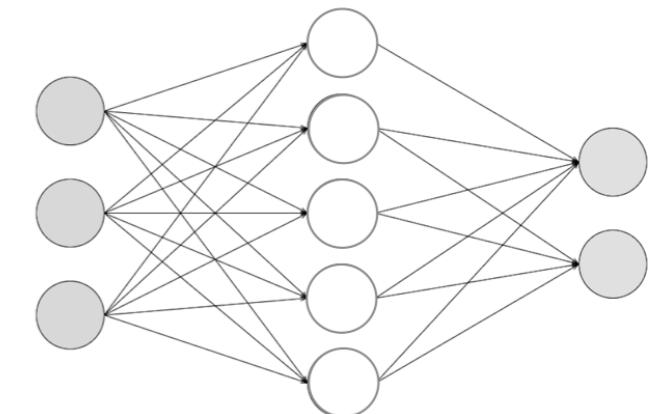
Sparse



Dense



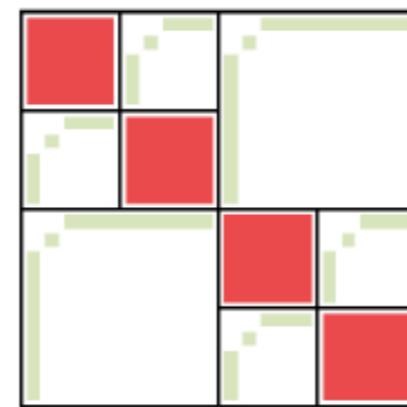
Million-D



locally connected



fully connected



group based on
connectivity

group based on
proximity

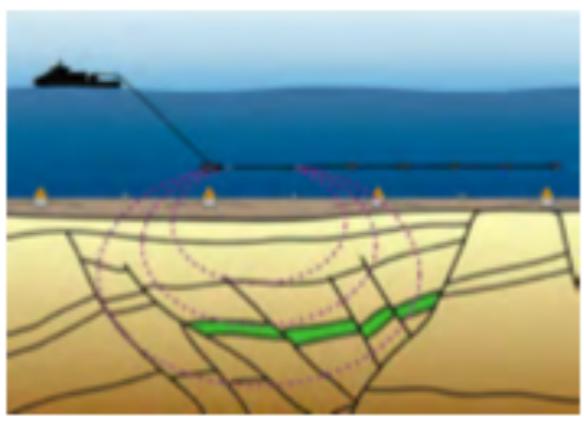
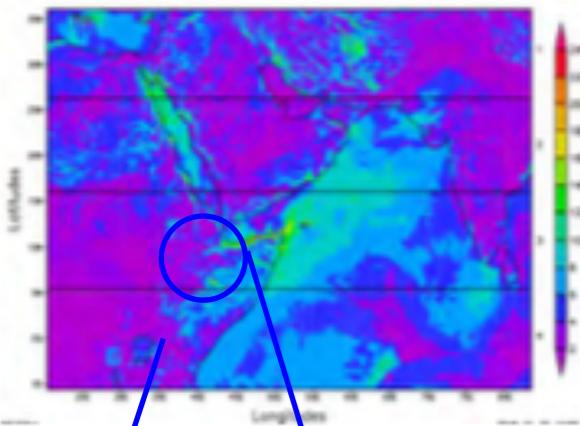
Each edge is a dimension

$$\begin{matrix} \text{large matrix} \\ = \end{matrix} \begin{matrix} \text{smaller matrix} \\ \otimes \\ \text{smaller matrix} \end{matrix}$$
A diagram showing a large matrix being factored into two smaller matrices using a tensor product symbol (\otimes).

everything is close but
dimensions are embedded

A different kind of piecewise linearity

Scientific computing



$$\int_{\Omega} f \phi d\Omega = 0$$

Conservation laws
are integrated over
low-D physical space

Local piecewise linearity

linear	linear	linear
linear	linear	linear
linear	linear	linear

$$Ax = b$$

Patch them together
and you get something
complex

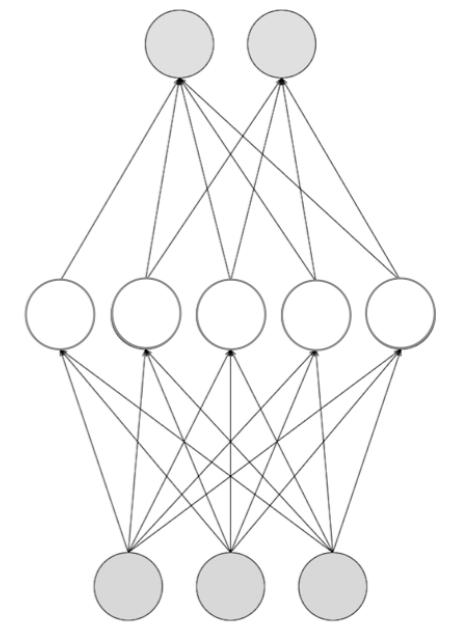
Deep learning

$$f(\boxed{\text{linear}})$$

$$g(\boxed{\text{linear}})$$

$$h(\boxed{\text{linear}})$$

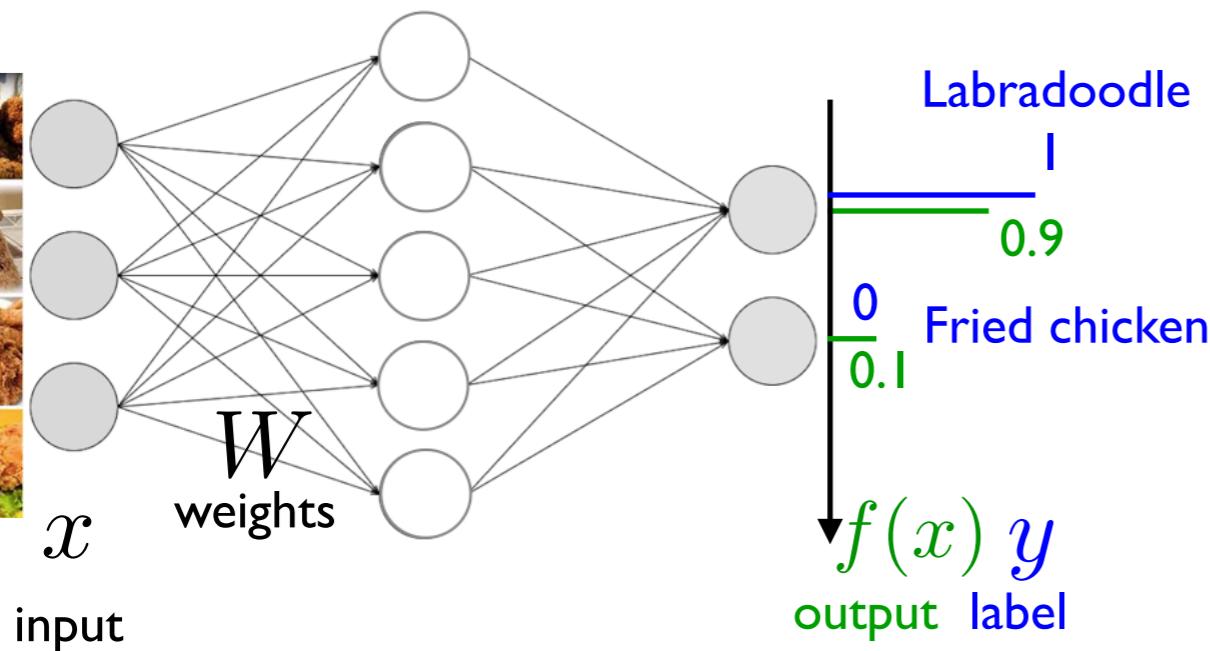
$$\mathbf{y} = f(g(h(\mathbf{x})))$$



Composite functions
of piecewise linear transformations
are used to describe complex
non-linear functions

Stack them up
and you get something
complex

Hessian, Fisher & Covariance Matrices



Negative log likelihood per class per data sample

$$l_{cd} = -\log f_c(x_d)$$

Loss per data sample

$$l_d = \sum_{class} y_{cl} l_{cd}$$

Overall loss

$$L = \sum_d l_d$$

Gradient

$$\nabla L = \sum_d \frac{\partial l_d}{\partial W}$$

Hessian (Newton's method)

$$H = \sum_d \frac{\partial^2 l_d}{\partial W^2}$$

Covariance (KFAC)

$$C = \sum_d \left(\frac{\partial l_d}{\partial W} \right)^T \left(\frac{\partial l_d}{\partial W} \right)$$

Fisher (Natural gradient descent)

$$F = \sum_d \sum_{class} f_c(x_d) \left(\frac{\partial l_{cd}}{\partial W} \right)^T \left(\frac{\partial l_{cd}}{\partial W} \right)$$

I'm intentionally forgetting vector notation to decouple calculus from linear algebra in this slide

Applications of H, F, C Matrices

Predicting Hyperparameters

An Empirical Model of Large-Batch Training

Sam McCandlish*
OpenAI
sam@openai.com

Jared Kaplan
Johns Hopkins University, OpenAI
jaredk@jhu.edu

Dario Amodei
OpenAI
damodei@openai.com

and the **OpenAI Dota Team**[†]

$$\mathcal{B}_{noise} = \frac{\text{tr}(\mathbf{H}\mathbf{C}^{-1})}{\mathbf{J}^T \mathbf{H} \mathbf{J}}$$

Optimizing Millions of Hyperparameters by Implicit Differentiation

Jonathan Lorraine Paul Vicol David Duvenaud
University of Toronto, Vector Institute
{lorraine, pvcoll, duvenaud}@cs.toronto.edu

$$\frac{\partial \theta}{\partial \lambda} = -\mathbf{H}^{-1} \frac{\partial^2 \mathcal{L}}{\partial \theta \partial \lambda^\top}$$

Preconditioned Optimizers

When Does Preconditioning Help or Hurt Generalization?

*Shun-ichi Amari[†], Jimmy Ba[‡], Roger Grosse[‡], Xuechen Li[§],
Atsushi Nitanda[¶], Taiji Suzuki[¶], Denny Wu[‡], Ji Xu^{||}

Gauss-Newton

$$\mathbf{F}(\theta)^{-1} \nabla \mathcal{L}(\theta) = \left\{ \mathbf{J}_{f,\theta}^\top \mathcal{H}_{\ell,f} \mathbf{J}_{f,\theta} \right\}^{-1} \mathbf{J}_{f,\theta}^\top \frac{\partial \mathcal{L}(\theta)}{\partial f}$$

Gram-Gauss-Newton

$$\mathbf{F}(\theta)^{-1} \nabla \mathcal{L}(\theta) = \mathbf{J}_{f,\theta}^\top \left\{ \mathcal{H}_{\ell,f} \mathbf{J}_{f,\theta} \mathbf{J}_{f,\theta}^\top \right\}^{-1} \frac{\partial \mathcal{L}(\theta)}{\partial f}$$

Bayesian Inference

Noisy Natural Gradient as Variational Inference

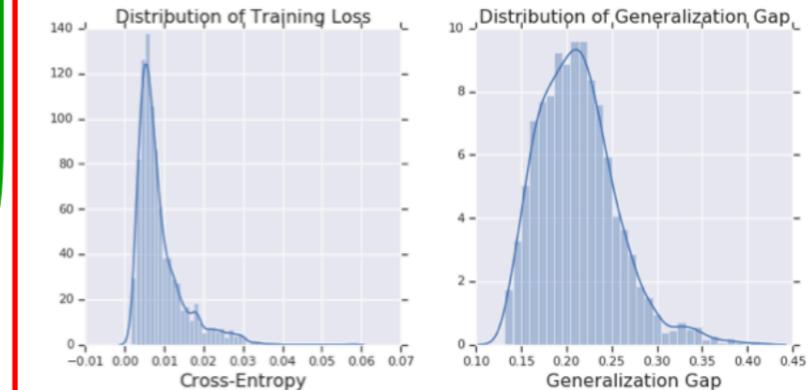
Guodong Zhang *^{1,2} Shengyang Sun *^{1,2} David Duvenaud ^{1,2} Roger Grosse ^{1,2}

$$\begin{aligned} & \left\{ \mathbf{F}(\theta) + \sigma^{-2} \mathbf{I} \right\}^{-1} \nabla \mathcal{L}(\theta) \\ &= \left\{ \mathbf{J}_{f,\theta}^\top \mathcal{H}_{\ell,f} \mathbf{J}_{f,\theta} + \sigma^{-2} \mathbf{I} \right\}^{-1} \mathbf{J}_{f,\theta}^\top \frac{\partial \mathcal{L}(\theta)}{\partial f} \end{aligned}$$

Generalization Metrics

*Fantastic Generalization Measures
and Where to find Them*

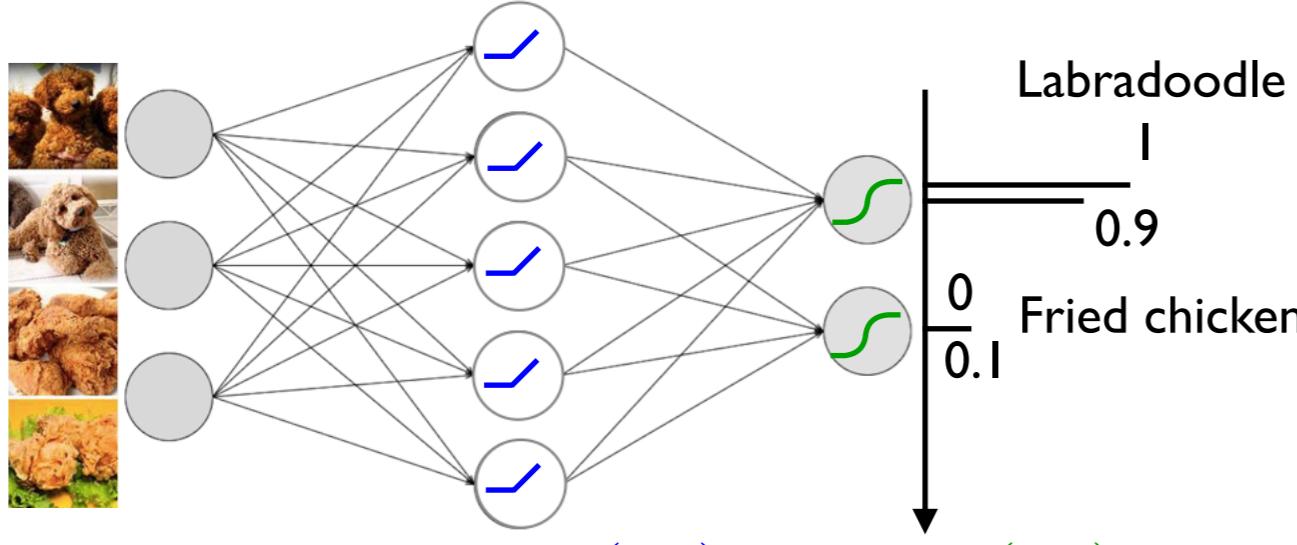
Yiding Jiang*, Behnam Neyshabur*, Hossein Mobahi
Dilip Krishnan, Samy Bengio



- Spectral bound
- Path norm
- Fisher-Rao metric
- Variance of gradients
- Sharpness
- PAC-Bayesian
- Takeuchi Information Criteria

$$\text{TIC}(\theta) = -\log p(y|\theta) + \frac{1}{N} \text{tr} (\mathbf{H}(\theta^*)^{-1} \mathbf{C}(\theta^*))$$

Jacobian-Vector Product



$$u_0 = W_0 h_0 \quad u_1 = W_1 h_1$$

$$\frac{\partial h_1}{\partial u_0} * \frac{\partial u_1}{\partial h_1} * \frac{\partial l_d}{\partial u_1} = \frac{\partial l_d}{\partial W}$$

$$\frac{\partial u_0}{\partial W_0} \quad \frac{\partial u_1}{\partial W_1}$$

2

$$\left(\frac{\partial u_1}{\partial W_0} \right)^T 25 * \begin{bmatrix} 1 \\ \frac{\partial l_d}{\partial u_1} \end{bmatrix}^2 = \begin{bmatrix} 1 \\ \frac{\partial l_d}{\partial W} \end{bmatrix}^{25}$$

Jacobian-vector product

Negative log likelihood per class per data sample

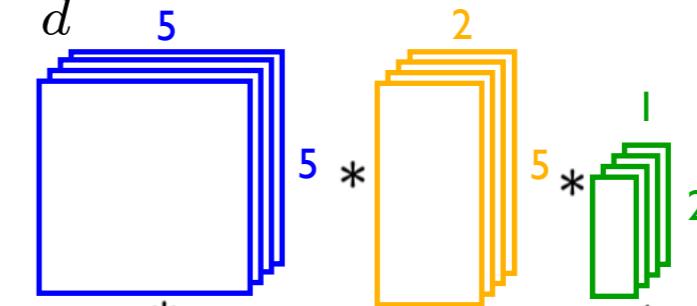
$$l_{cd} = -\log f_c(x_d)$$

Loss per data sample

$$l_d = \sum_{class} y_{cl} l_{cd}$$

Overall loss

$$L = \sum_{data} l_d$$



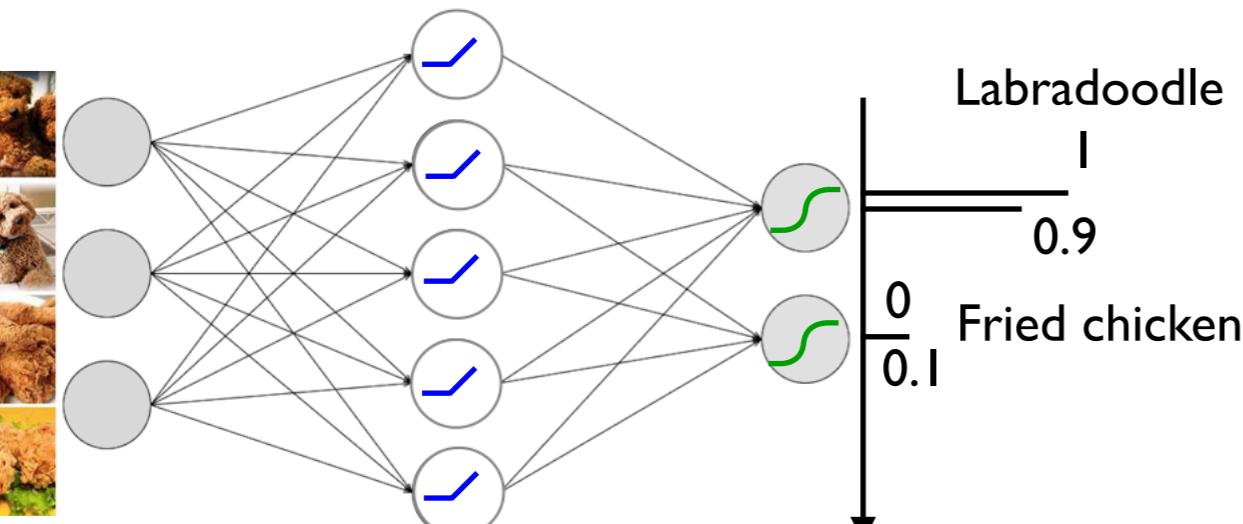
$$\sum_{data} l_d$$

Gradient

$$\nabla L = \sum_{data} \frac{\partial l_d}{\partial W}$$

$$\sum_{data} \frac{\partial l_d}{\partial W}$$

Jacobian-Matrix Product



$$x = h_0 \quad h_1 = f(u_0) \quad p = f(u_1)$$

$$u_0 = W_0 h_0 \quad u_1 = W_1 h_1$$

Gauss-Newton approximation →

$$\left(\frac{\partial u_1}{\partial W_0} * \sqrt{\frac{\partial^2 l_d}{\partial u_1^2}} \right)^T \left(\sqrt{\frac{\partial^2 l_d}{\partial u_1^2}} * \frac{\partial u_1}{\partial W_0} \right)$$

$$\left(\frac{\partial u_1}{\partial W_0} \right)^T 25 * \boxed{\sqrt{\frac{\partial^2 l_d}{\partial u_1^2}}}^2$$

Jacobian-matrix product

Gradient of first layer per data sample

$$\frac{\partial l_d}{\partial W_0} = \frac{\partial u_0}{\partial W_0} * \frac{\partial h_1}{\partial u_0} * \frac{\partial u_1}{\partial h_1} * \frac{\partial l_d}{\partial u_1}$$

Hessian of first layer per data sample

$$\frac{\partial^2 l_d}{\partial W_0^2} = \frac{\partial^2 u_0}{\partial W_0^2} * \frac{\partial h_1}{\partial u_0} * \frac{\partial u_1}{\partial h_1} * \frac{\partial l_d}{\partial u_1} \rightarrow 0$$

$$+ \left(\frac{\partial u_0}{\partial W_0} \right)^2 * \frac{\partial^2 h_1}{\partial u_0^2} * \frac{\partial u_1}{\partial h_1} * \frac{\partial l_d}{\partial u_1} \rightarrow 0$$

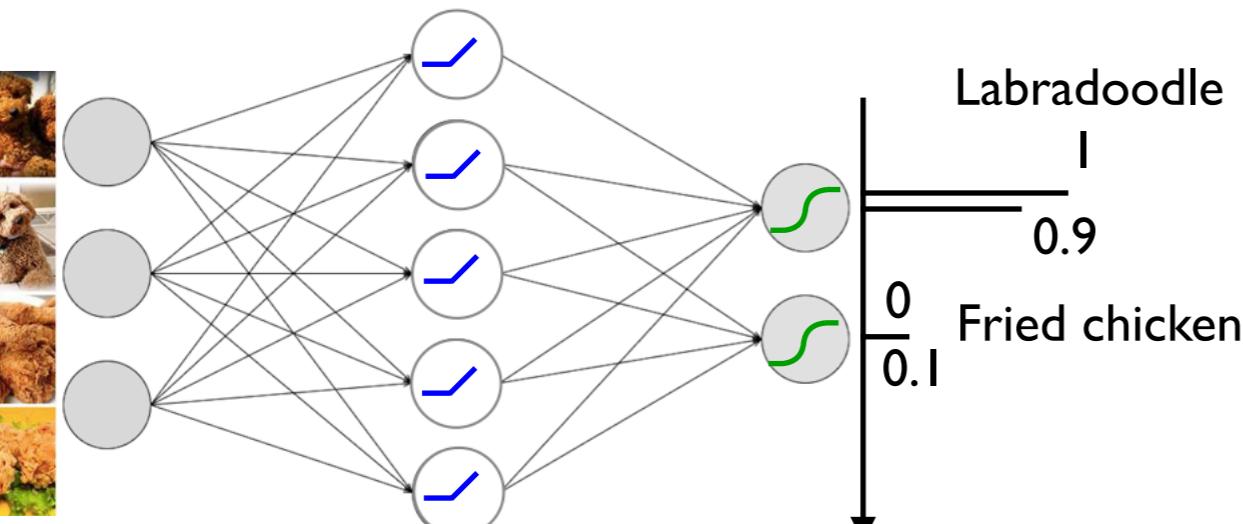
$$+ \left(\frac{\partial u_0}{\partial W_0} \right)^2 * \left(\frac{\partial h_1}{\partial u_0} \right)^2 * \frac{\partial^2 u_1}{\partial h_1^2} * \frac{\partial l_d}{\partial u_1} \rightarrow 0$$

$$+ \left(\frac{\partial u_0}{\partial W_0} \right)^2 * \left(\frac{\partial h_1}{\partial u_0} \right)^2 * \left(\frac{\partial u_1}{\partial h_1} \right)^2 * \frac{\partial^2 l_d}{\partial u_1^2}$$

$$= \left(\frac{\partial u_1}{\partial W_0} \right)^T * \frac{\partial^2 l_d}{\partial u_1^2} * \left(\frac{\partial u_1}{\partial W_0} \right)$$

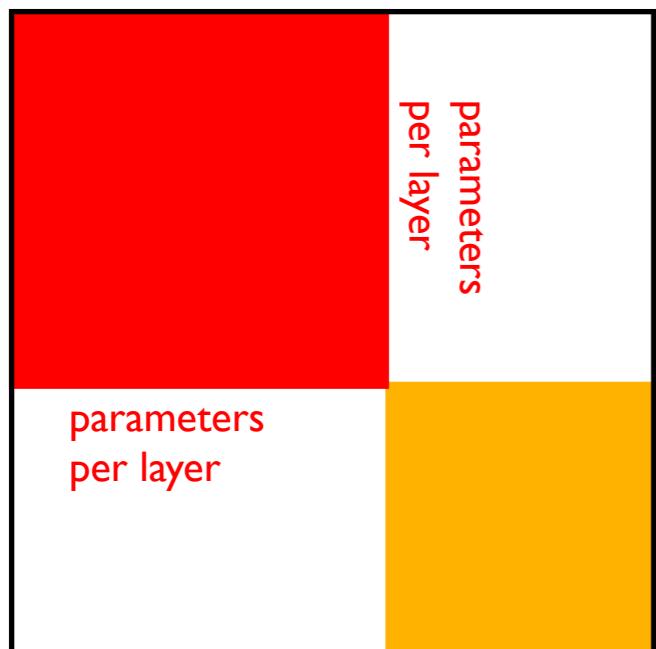
$$\left(\frac{\partial u_1}{\partial W_0} \right)^T 25 * \boxed{\frac{\partial^2 l_d}{\partial u_1^2}}^2 * \left(\frac{\partial u_1}{\partial W_0} \right)^2$$

Kronecker Factors

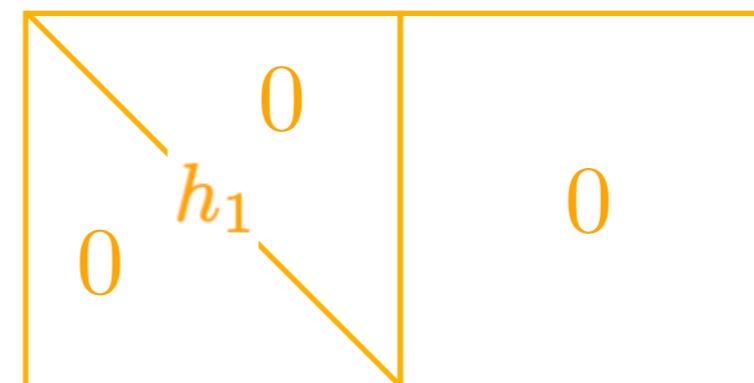


$$x = h_0 \quad h_1 = f(u_0) \quad p = f(u_1)$$

$$u_0 = W_0 h_0 \quad u_1 = W_1 h_1$$



$$\begin{pmatrix} \frac{\partial u_1}{\partial W_1} & \frac{\partial u_1}{\partial W_2} & \cdots & \frac{\partial u_1}{\partial W_{10}} \\ \frac{\partial u_2}{\partial W_1} & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ \frac{\partial u_5}{\partial W_1} & \cdots & \cdots & \frac{\partial u_5}{\partial W_{10}} \end{pmatrix} \in \mathbb{R}^{10 \times 5}$$



$$\begin{aligned} \frac{\partial l_d}{\partial W_0} &= \frac{\partial u_0}{\partial W_0} \otimes \frac{\partial h_1}{\partial u_0} * \frac{\partial u_1}{\partial h_1} * \frac{\partial l_d}{\partial u_1} \\ &= \frac{\partial u_0}{\partial W_0} \otimes \frac{\partial l_d}{\partial u_0} \\ \frac{\partial l_d}{\partial W_1} &= \frac{\partial u_1}{\partial W_1} \otimes \frac{\partial l_d}{\partial u_1} \end{aligned}$$

$$\begin{matrix} \text{Large Matrix} \\ = \end{matrix} \begin{matrix} \text{Matrix 1} \\ \otimes \\ \text{Matrix 2} \end{matrix}$$

K-FAC

Kronecker-Factored Approximate Curvature (J. Martens+, ICML 2015)

Step 1. Layer-wise block-diagonal approximation

$$C = F \approx \text{FIM}$$

Step 2. Kronecker-factorization (for each layer)

$$F_\ell \approx \square \otimes \square$$

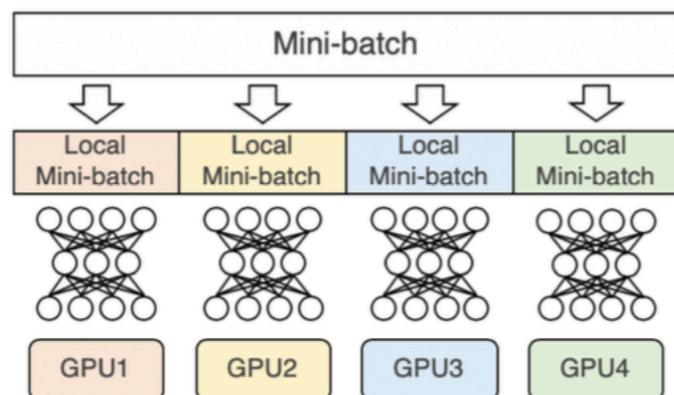
Kronecker product

$$\left(\begin{array}{c|c} \square & \square \\ \hline \square & \square \end{array} \otimes \begin{array}{c|c} \square & \square \\ \hline \square & \square \end{array} \right)^{-1} = \begin{array}{c|c} \square & \square \\ \hline \square & \square \end{array} \otimes \begin{array}{c|c} \square & \square \\ \hline \square & \square \end{array}^{-1}$$

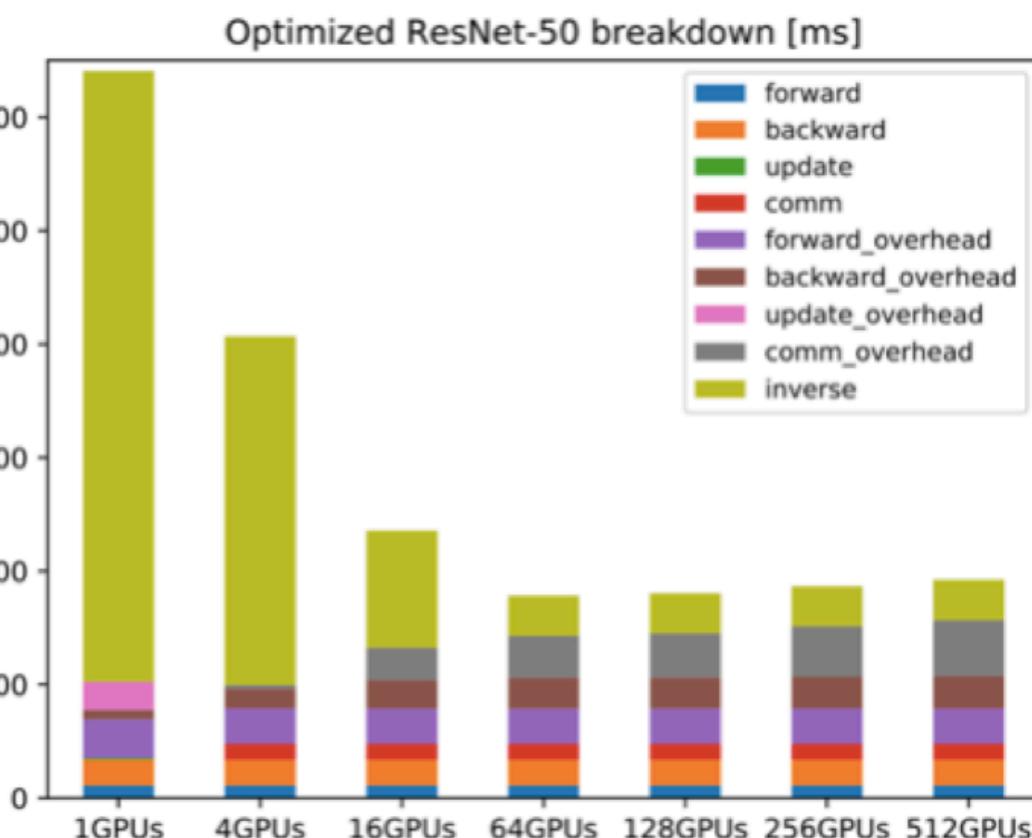
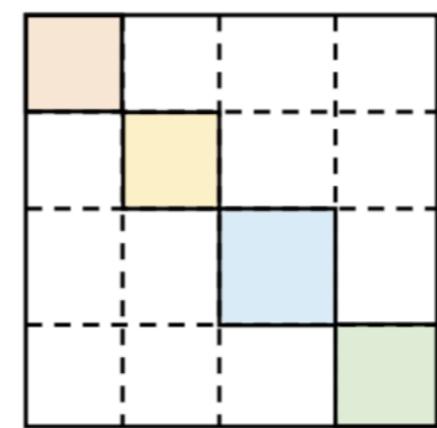
much easier to invert $O(N_\ell^3) \rightarrow O(N_\ell^{3/2})$

Distributed K-FAC (K. Osawa et al., CVPR2019, TPAMI)

data-parallelism \times layers-parallelism



$$F \approx$$



Summary

- The structure of matrices depend on the underlying geometry
- Sparsity is related to connectivity, where as rank is related to proximity
- FMM is a matrix-free mat-vac of a hierarchical low-rank matrix
- Using the matrix form is useful for multiple right hand sides and factorization
- Nesting of bases and admissibility distinguishes the various hierarchical formats
- FMM and hierarchical matrices both rely on geometric separation
- Everything is close in high dimensions but dimensions may also have structure
- Embedding of dimensions leads to Kronecker structure
- Exploiting this structure is useful for very high-dimension problems

Thank you

