

第40回 計算数理工学フォーラム

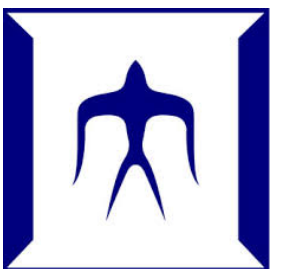
2021/9/24

# 階層的低ランク近似法に関するレビュー

東京工業大学 学術国際情報センター

横田 理央

[rioyokota@gsic.titech.ac.jp](mailto:rioyokota@gsic.titech.ac.jp)



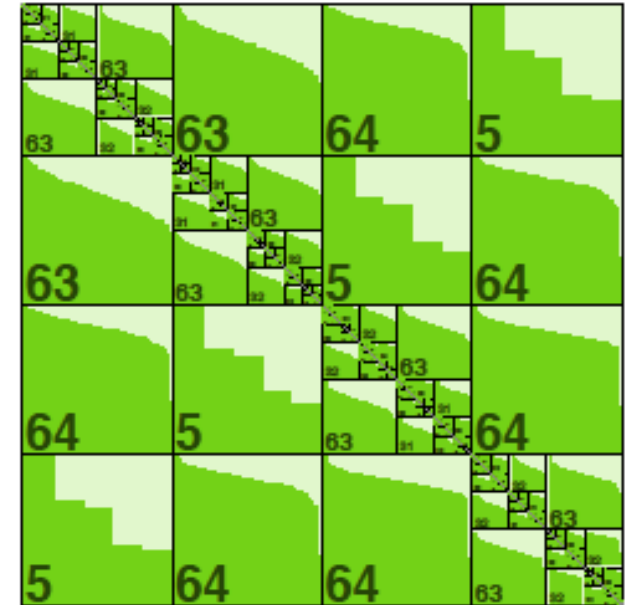
# 階層的低ランク近似法とは？

## 密行列の近似直接解法

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

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

非対角ブロックのランクは小さいという仮定



## 疎行列に使えるのか？

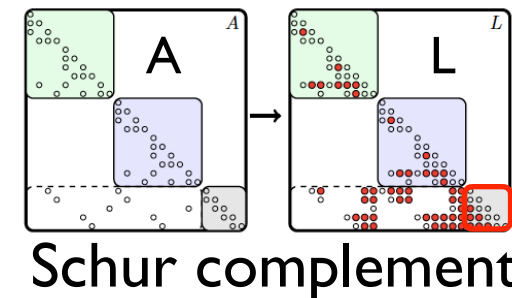
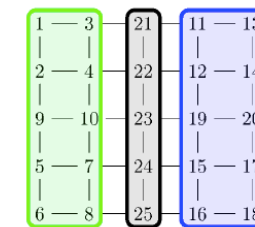
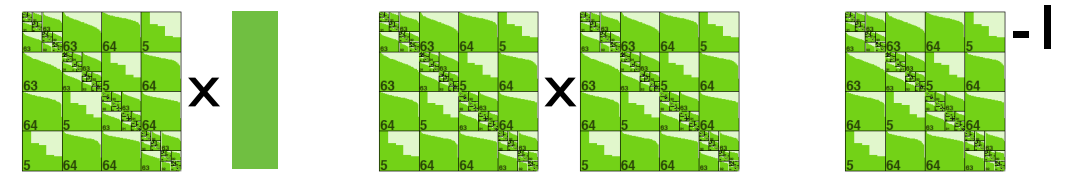
Fill-inさせないことの方が重要

Schur補元は密だが非対角ブロックのランクは小さい

反復法では？

前処理として使えるがMultigridにはなかなか勝てない

条件数が悪い問題でなら優位かも？

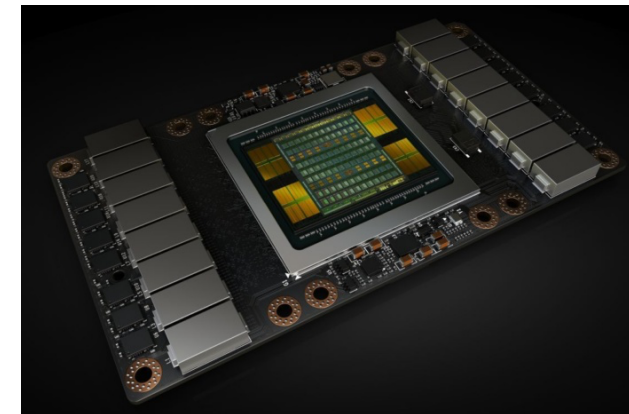
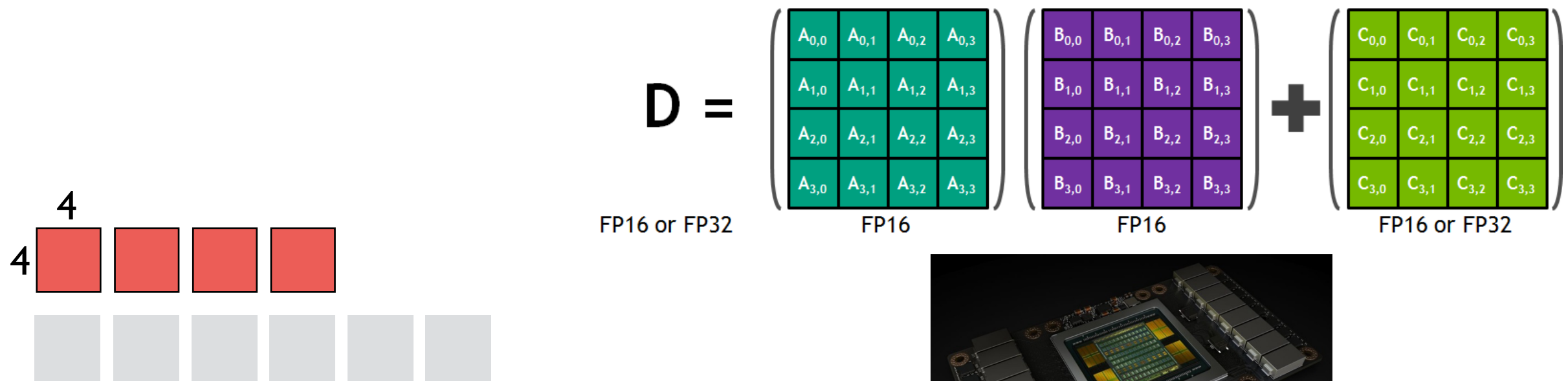


# 機械学習向けハードウェアと相性が良い

## 小さな密行列演算がたくさん生じる

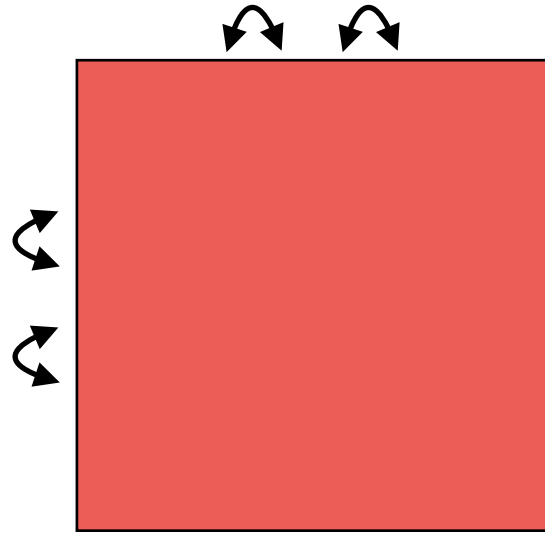


## 低ランク近似なので低精度演算で十分



# 階層的低ランク近似法の3つのステップ

## 並べ替え



最小化したいのは？

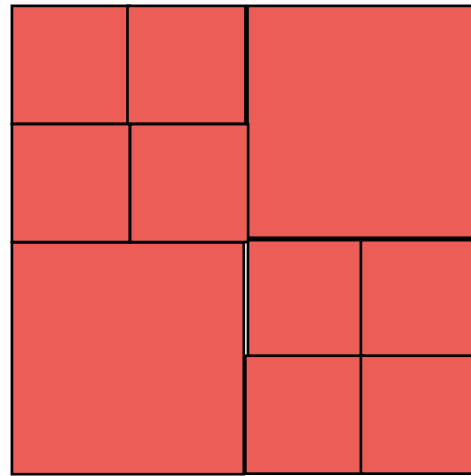
ランク (幾何学的距離)

通信 (データ局在性)

Fill-in (グラフの接続)

→ 普通は近いものと繋がっているので Fill-in とランクの最小化は両立する

## 階層化



どこまで分割するか？

分割数を増やせばそれぞれのランクは小さくなる

ランクは固定で分割によって精度を制御することもできる

→ SIMD friendly

## 低ランク近似



速さか, 安定性か？

ACAは速いが不安定  
RSVDは安定だが遅い

どのブロックを近似するかを選択するのに Gram 距離を使うこともできる

# Replacing Exact Linear Algebra with Low-Rank

厳密解

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

近似解

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

Application

App.

ScaLAPACK

cuSolverMG

LAPACK  
PLASMA

MKL

cuSolverDN  
MAGMA

BLAS

CUBLAS

CPU

FP64

GPU

FP32



分散

QR

LU

MatMul

Mat-vec

STRUMPACK

HiCMA

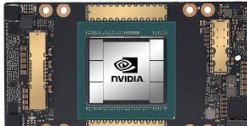
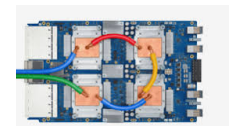
GOFMM

LoRaSp

HBLAS

?PU

TF32, bfloat16



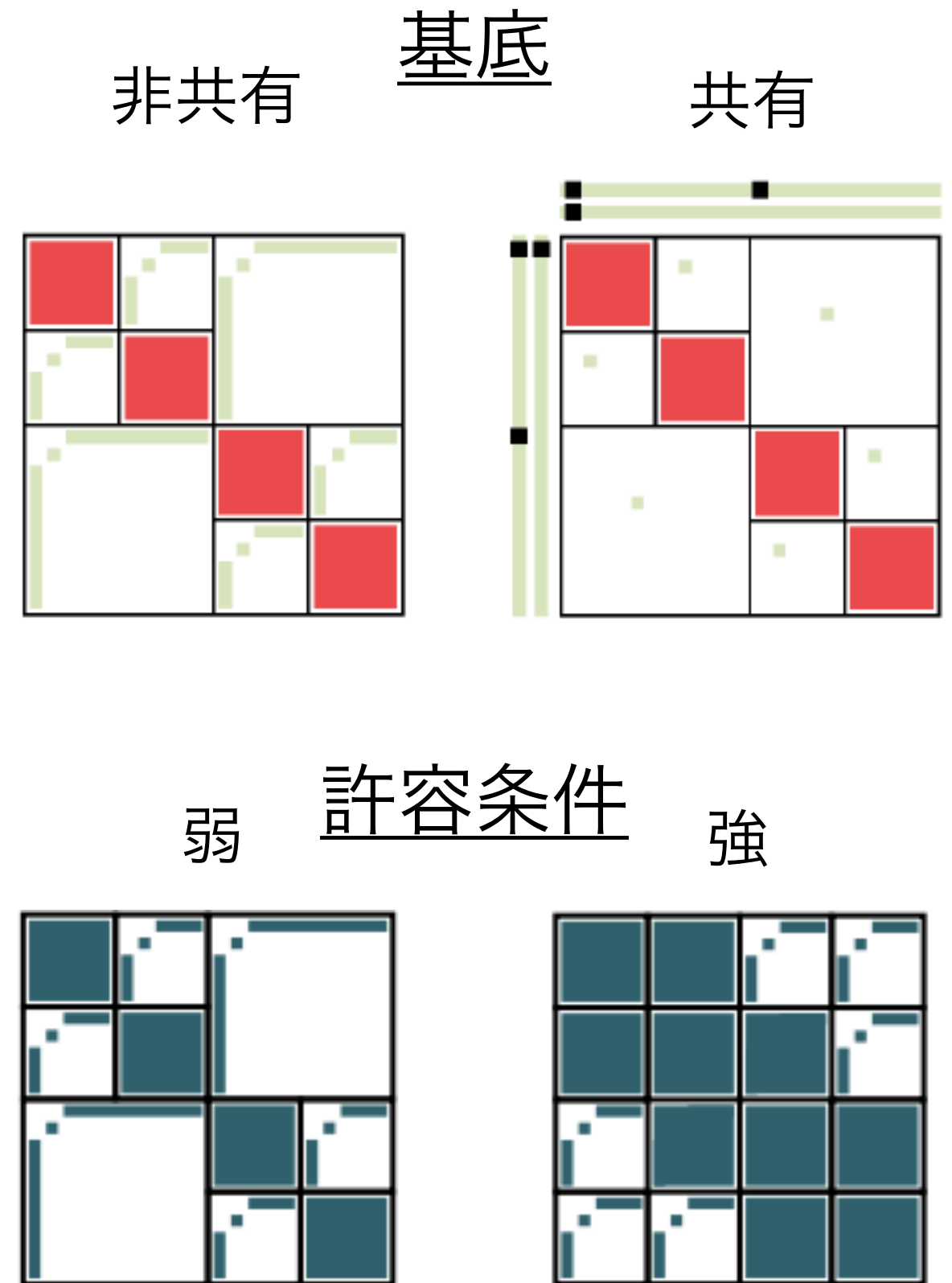
# List of implementations

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

[https://github.com/gchavez2/awesome\\_hierarchical\\_matrices](https://github.com/gchavez2/awesome_hierarchical_matrices)

# 手法間の違い

	基底の 共有	許容条件
H行列	無	強
H <sup>2</sup> 行列	有	強
HODLR	無	弱
HSS	有	弱
BLR	無	非階層的
BLR <sup>2</sup>	有	非階層的

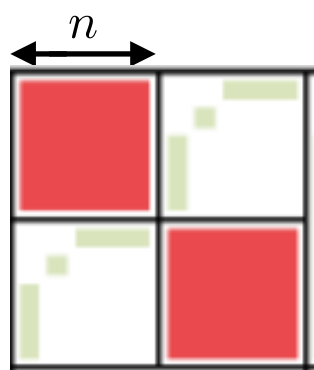


# Nullity Theorem

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

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

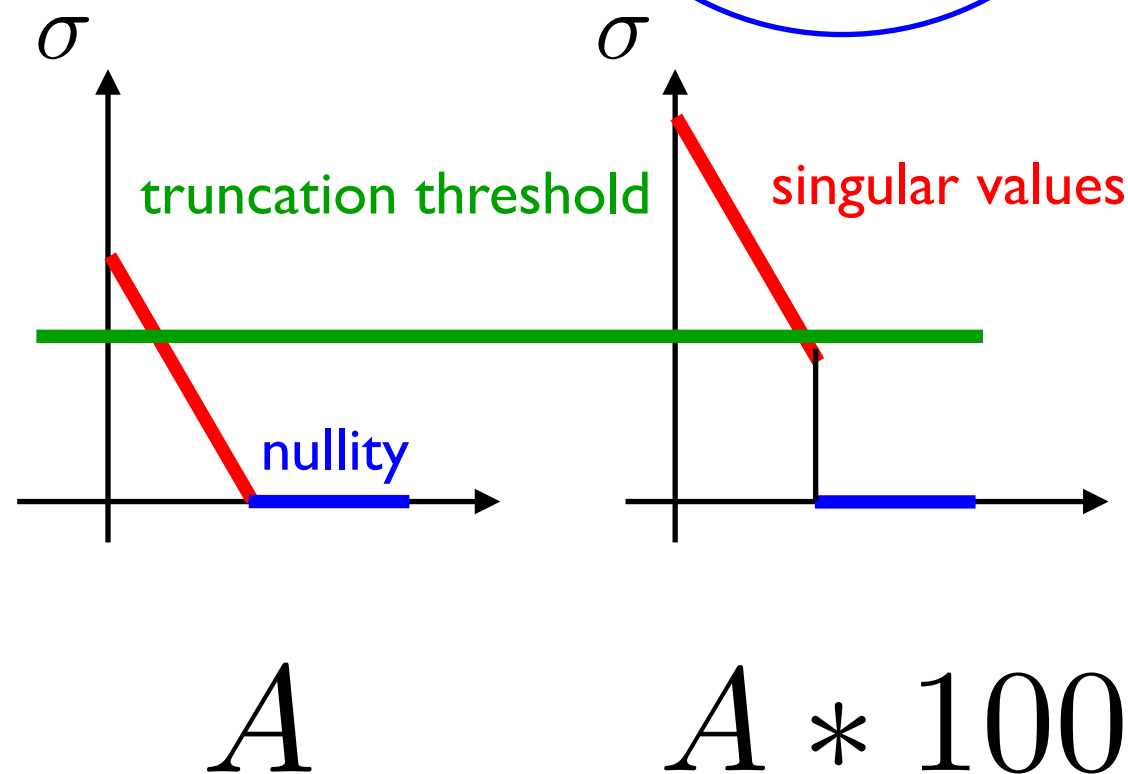
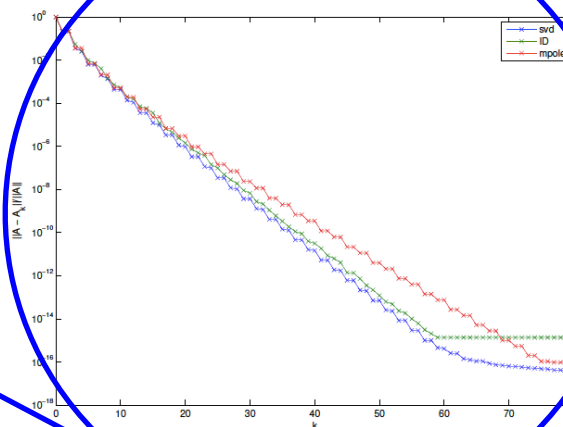
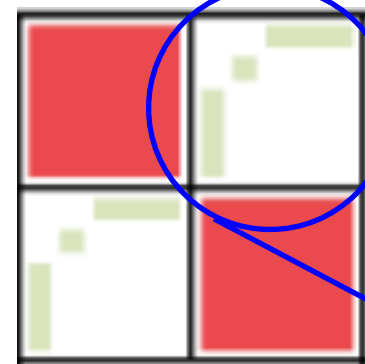
$$\text{rank}(A) + \text{nullity}(A) = n.$$



Apply it recursively

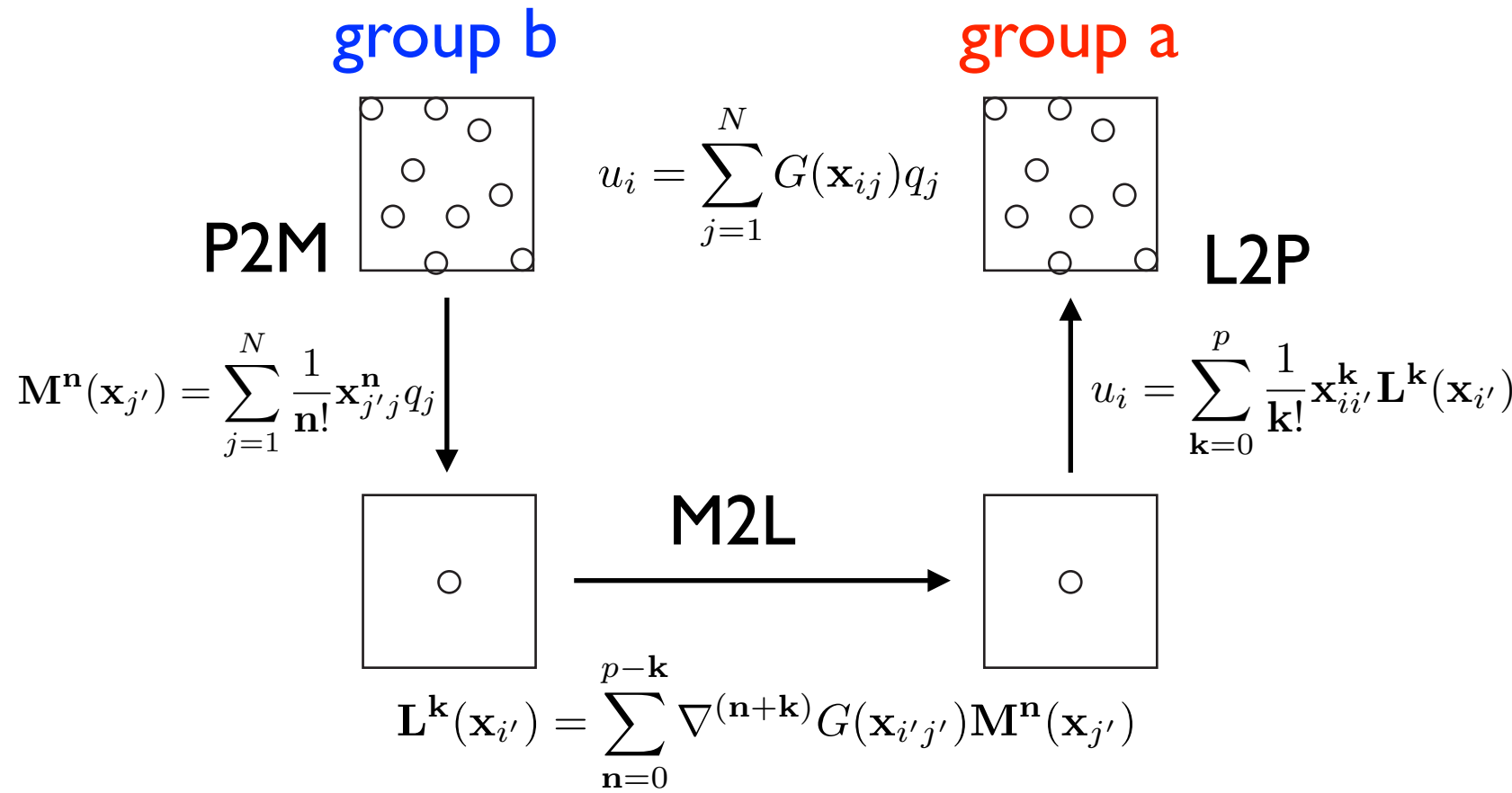
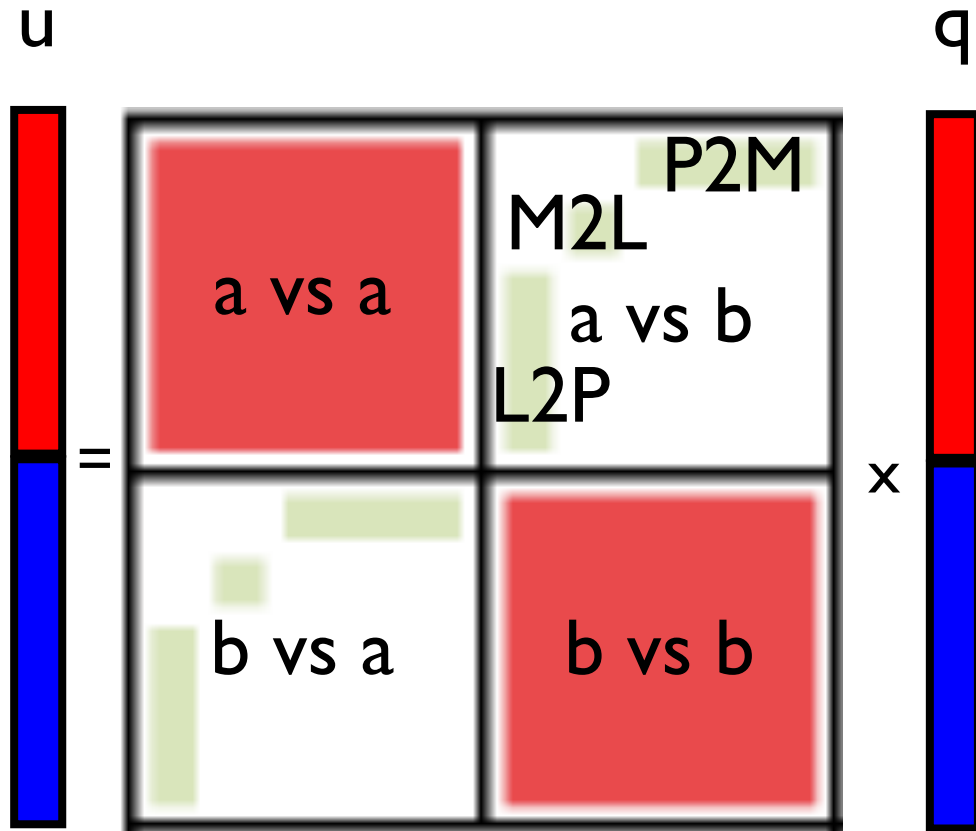


特異値の減衰

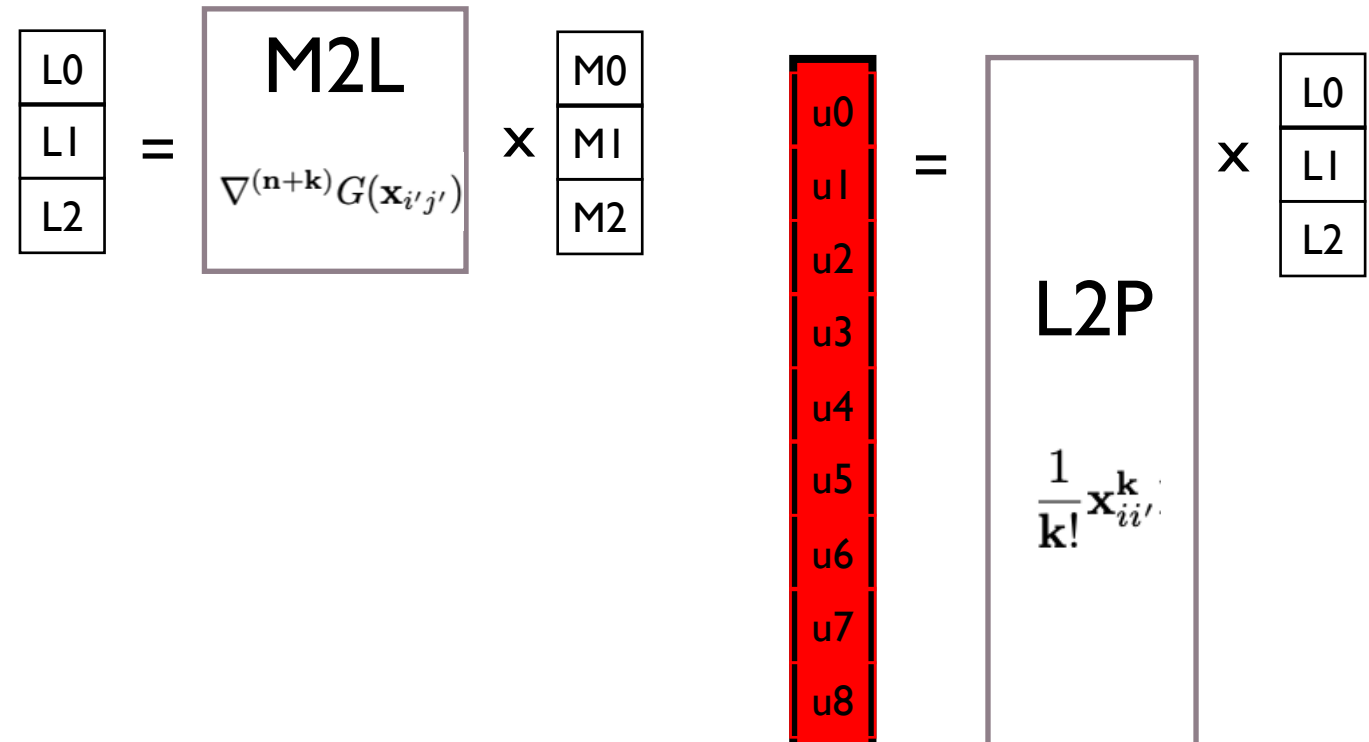
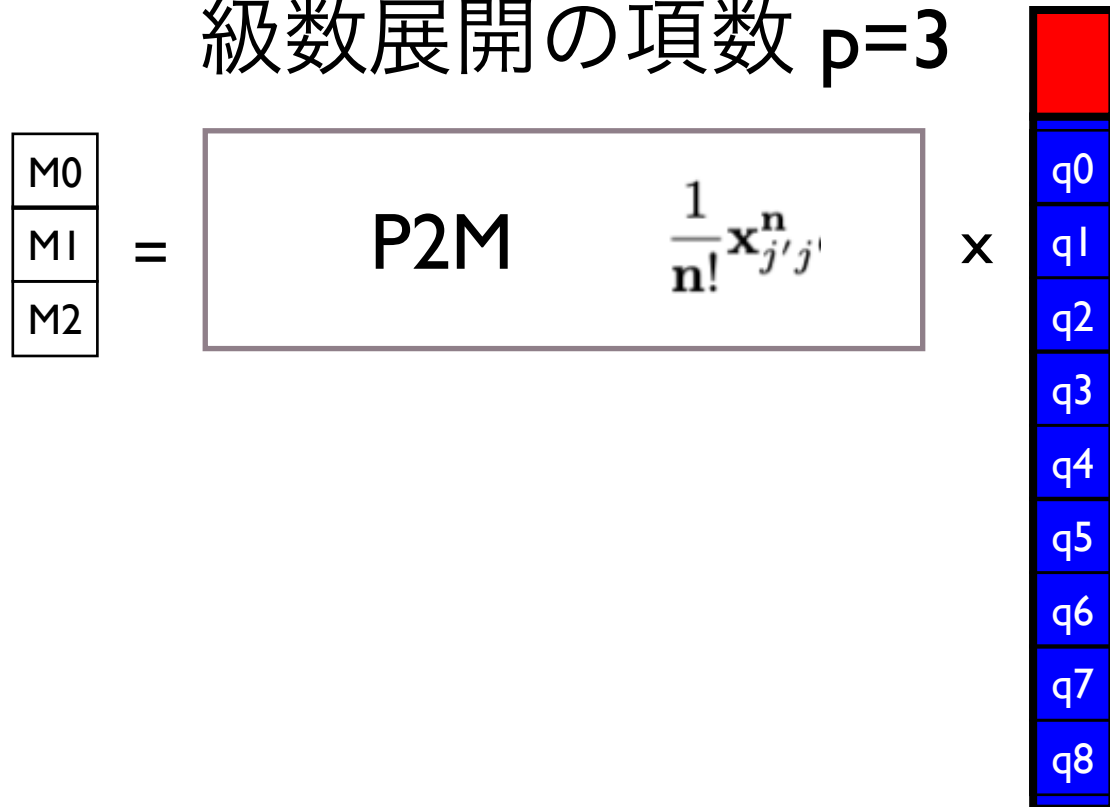




# FMMはH<sup>2</sup>行列-ベクトル積

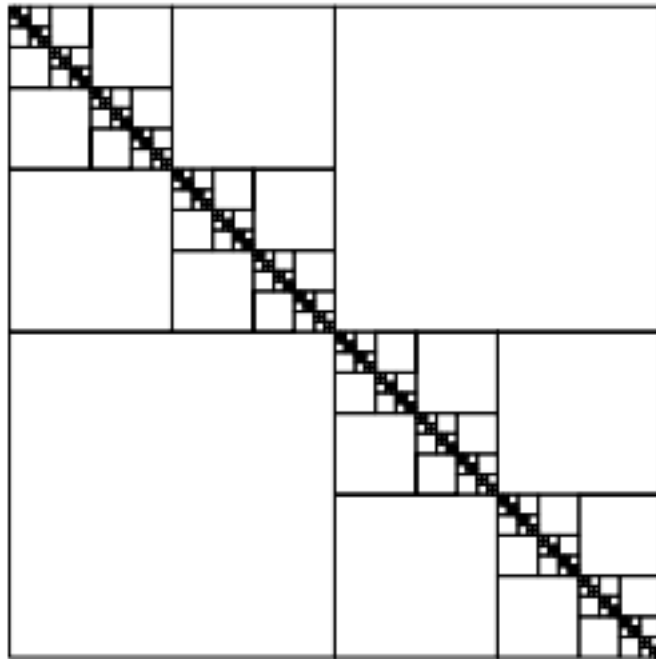


級数展開の項数  $p=3$

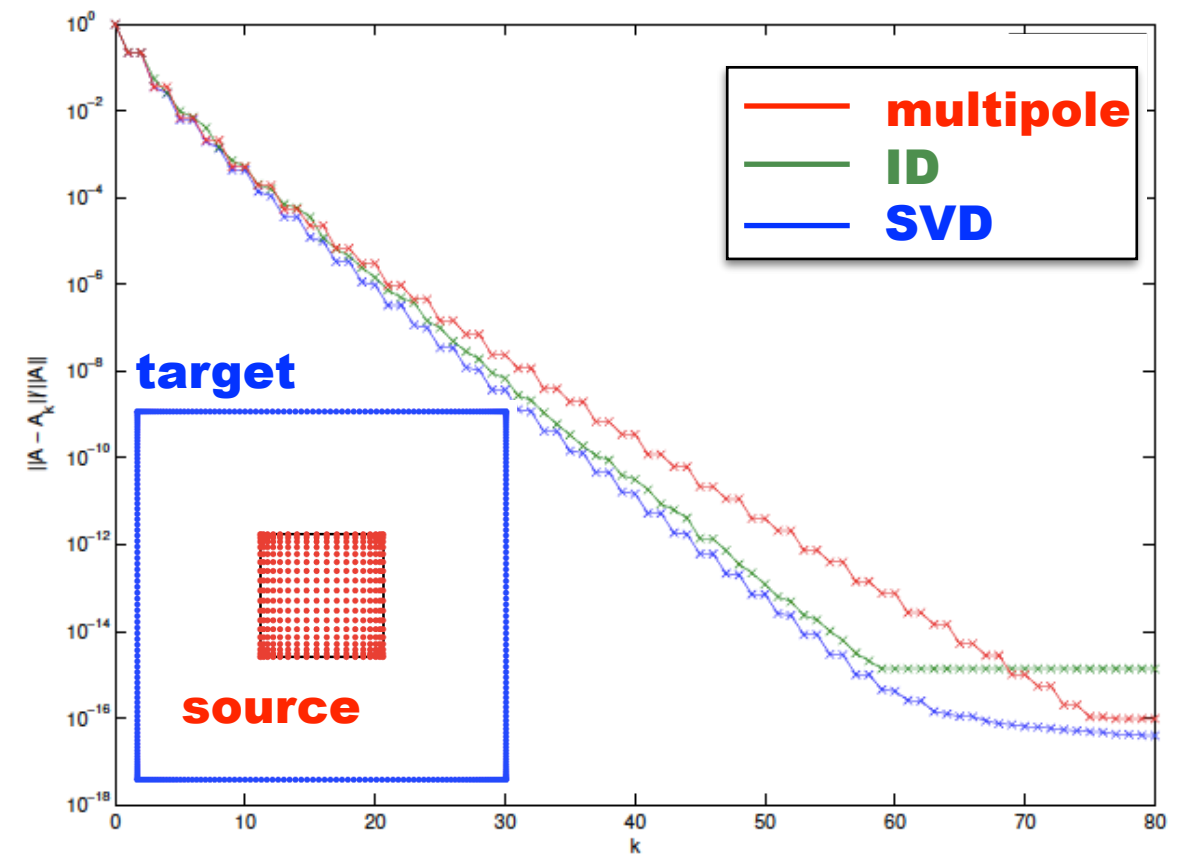
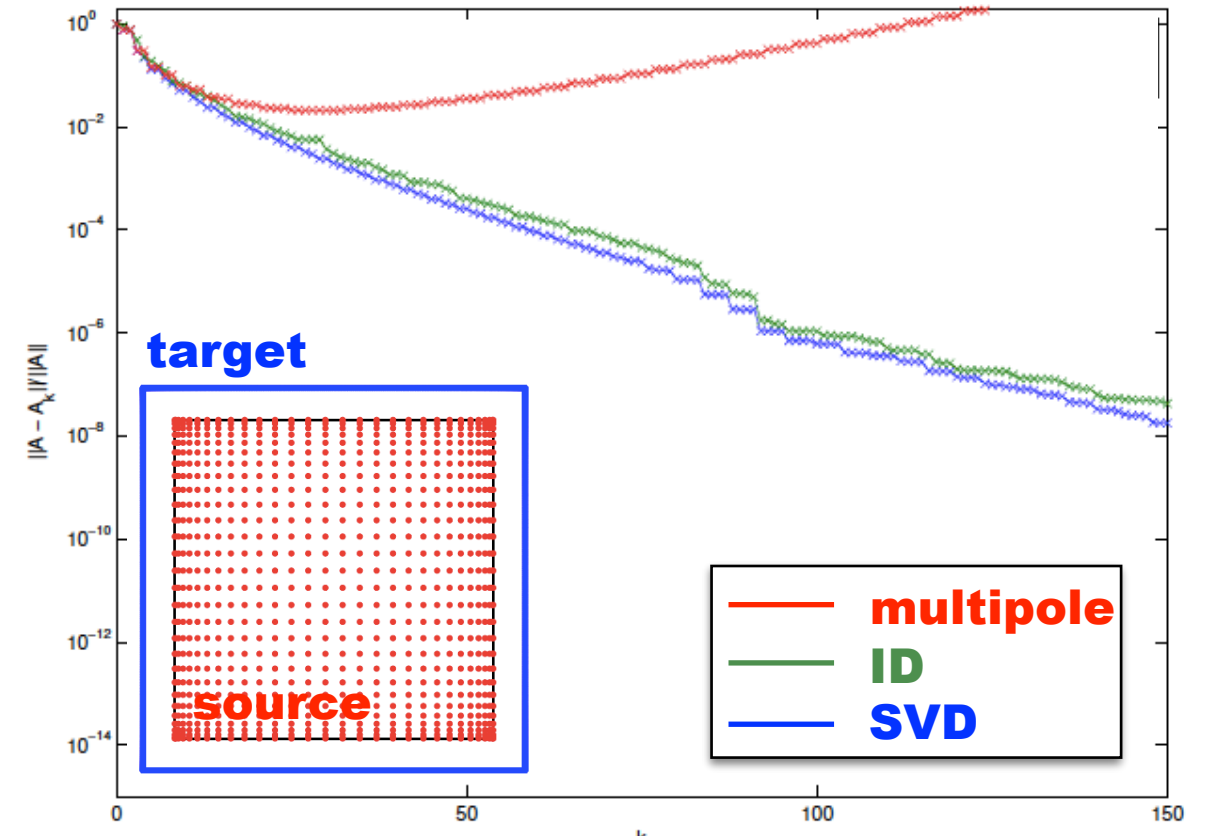
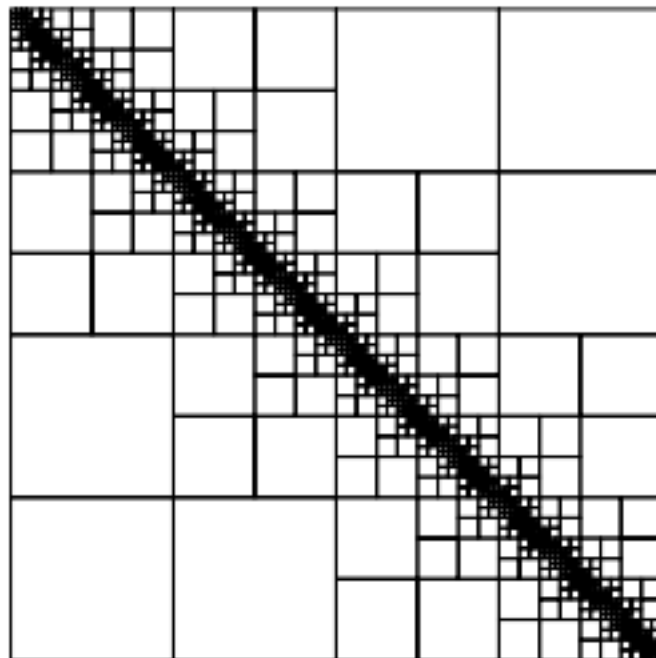


# 許容条件

## Weak admissibility

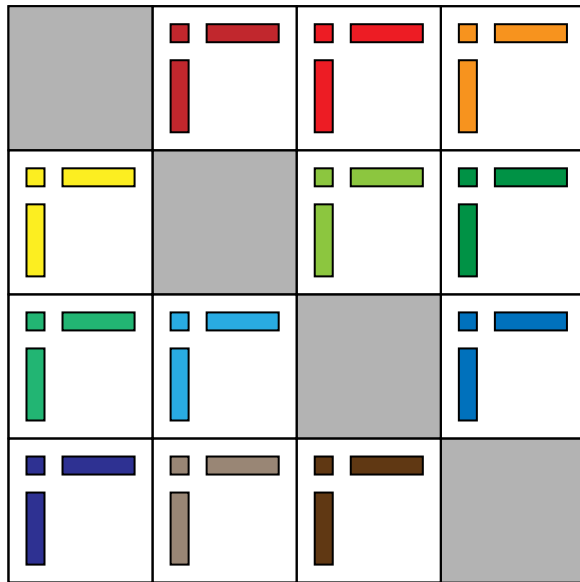


## Standard admissibility

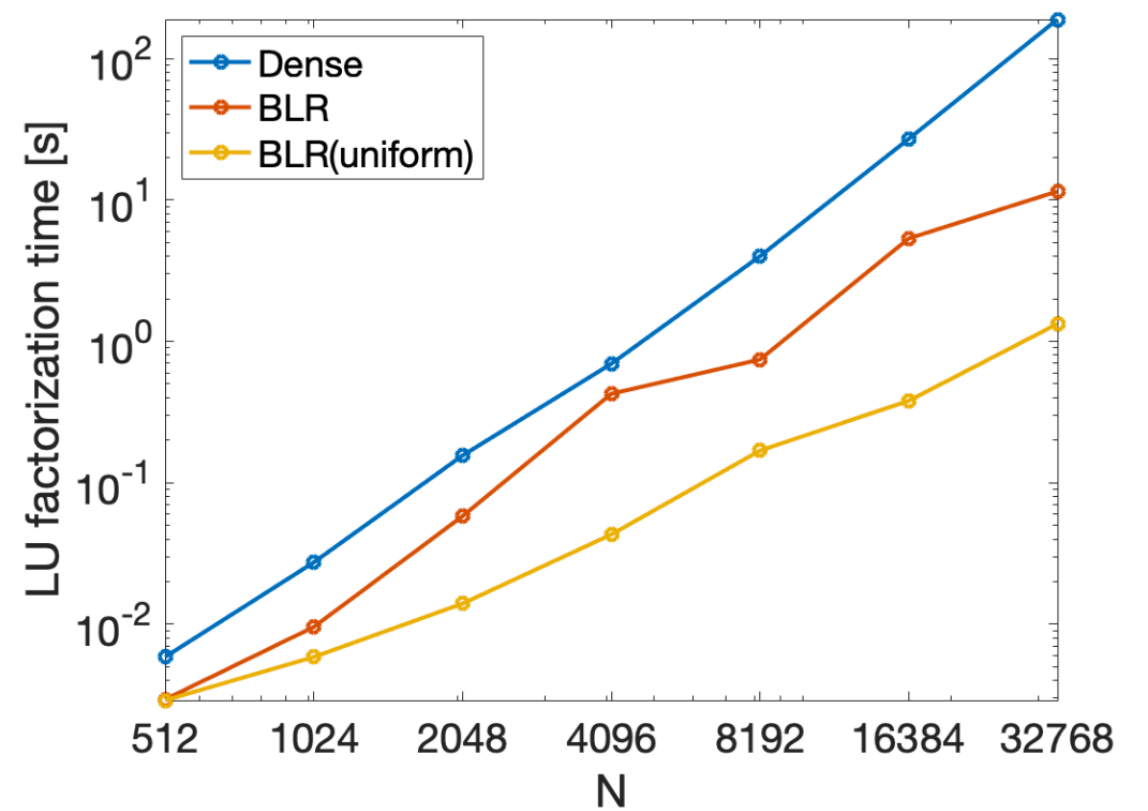
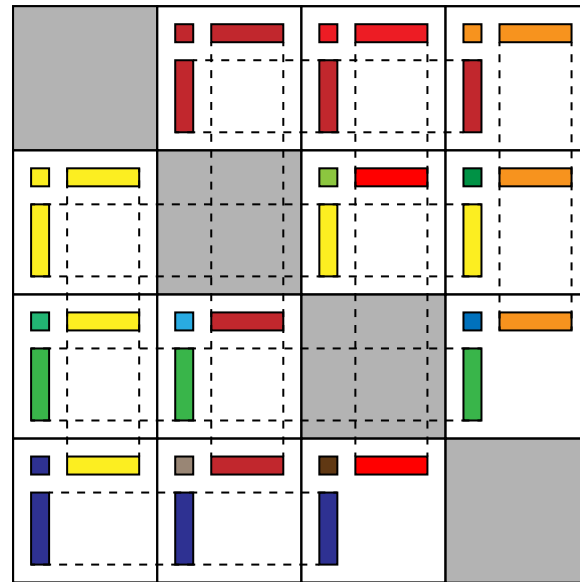


# Uniform Basis

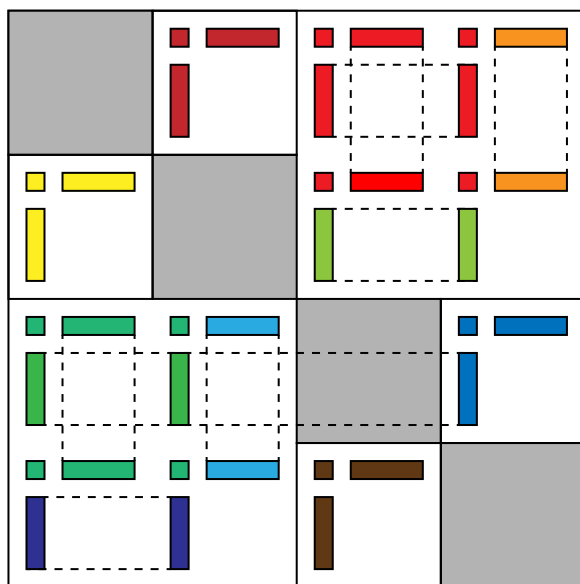
BLR



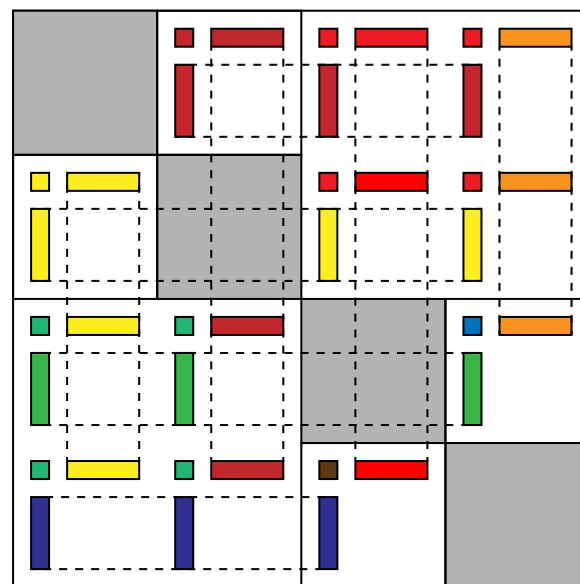
BLR<sup>2</sup>



H-matrix (HODLR)

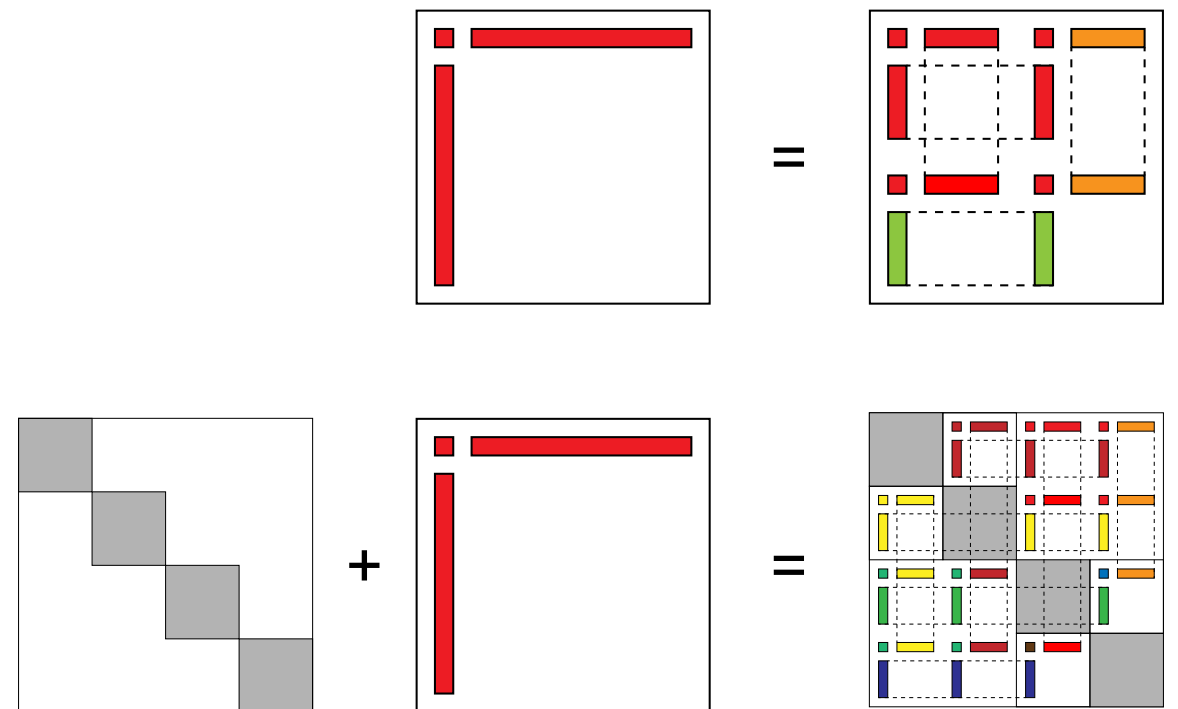


H<sup>2</sup>-matrix (HSS)



Nonuniform basis

Uniform basis



# 勢力図

## Germany

Shared memory H-LU  
Kriemann (2014)

Nested cross approximation  
Börm & Christophersen (2014)

H<sup>2</sup>-matrix for eigenvalues  
Berner et al. (2015)

OmpSs H-LU  
Aliaga et al. (2017)

GCA H<sup>2</sup>-matrix  
Börm et al. (2018)

## Berkeley

HSS2D  
Xia (2014)

HSS selected inversion  
Xia et al. (2015)

Superfast DC eigenvalue  
Vogel, et al. (2016)

Shared memory HSS MF  
Ghysels et al. (2016)

Distributed HSS MF  
Rouet et al. (2016)

## Japan

Distributed H-matrix  
Ida et al. (2015)

Distributed GPU H-matrix  
Yamazaki et al. (2018)

Lattice H-matrix  
Ida (2018)

GPU load-balance ACA  
Hoshino et al. (2018)

Mixed precision H-matrix  
Ooi et al. (2020)

## EPFL

HODLR QR  
Kressner et al. (2018)

## Minnesota

Multilevel Low-Rank  
Li & Saad (2013)

DD Low-Rank  
Li & Saad (2014)

Multilevel Schur Low-Rank  
Xi et al. (2016)

SMASH  
Cai et al. (2018)

H<sup>2</sup>-matrix + FMM  
Xing & Chow (2021)

## Stanford(Ying)

O(N) RS 2-D  
Corona (2015)

HIF for PDEs  
Ho & Ying (2016)

Distributed memory HIF  
Li & Ying (2016)

RS for maximum likelihood  
Minden et al. (2016)

RS with strong admissibility  
Minden et al. (2017)

Quantized Tensor Train  
Corona et al. (2017)

## Stanford(Darve)

HODLR multifrontal  
Aminfar et al. (2016)

IFMM preconditioned Stokes  
Coulter et al. (2017)

IFMM preconditioned Helmholtz  
Takahashi et al. (2017)

Non-extensive sparsification  
Sushnikova et al. (2017)

Sparsified Nested Dissection  
Cambier et al. (2019)

spaND QR  
Gnanasekaran et al. (2020)

## INRIA

BLR multifrontal  
Amestoy et al. (2015)

BLR multicore  
Amestoy et al. (2017)

Multilevel BLR  
Amestoy et al. (2019)

## KAUST

BLR Cholesky  
Akbulduk et al. (2017)

Batched QR, SVD  
Boukaram et al. (2018)

GPU MatVec  
Boukaram et al. (2019)

## Texas (Biros)

inv-ASKIT  
Yu et al. (2016)

Distributed inv-ASKIT  
Yu et al. (2017)

GOFMM  
Yu et al. (2017)

Distributed GOFMM  
Yu et al. (2018)

# 最近の論文の紹介

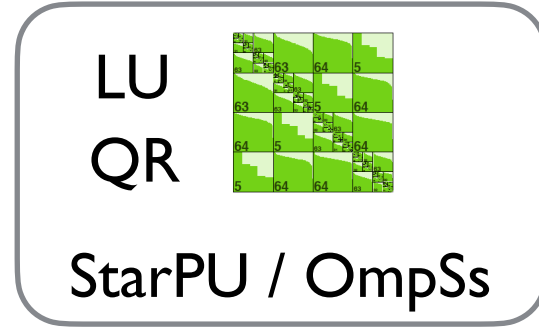
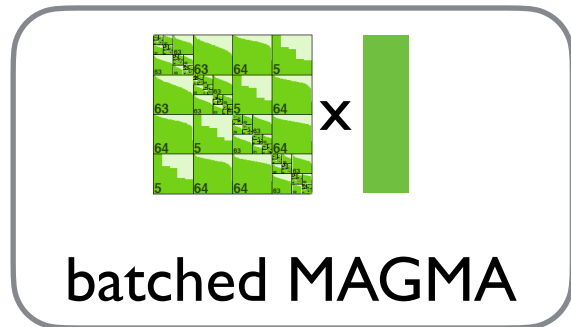
著者・年	行列	演算	構造	近似法	OpenMP	MPI	GPU
2016Vogel	密	QAQT	HSS	?	×	×	×
2017Akbulduk	密	Cholesky	BLR	RSVD	○	×	×
2017Fernando	密	LU	HSS	ID	×	○	○
2017Ghysels	疎	LU	HSS	RSVD	○	○	×
2017Li	疎	LU	HSS	RRQR	×	○	×
2017Minden	疎	LU	H2	RRQR	×	×	×
2018Amestoy	疎	LU	MBLR	?	×	×	×
2018Börm	密	O(N)圧縮	H2	GCA	×	×	×
2018Cai	密	O(N)圧縮	HSS	RRQR	×	×	×
2018Kressner	密	QR	HODLR	?	×	×	×
2018Yu	密	O(N)圧縮	HSS	ID	○	○	×
2019Amestoy	疎	LU	BLR	?	○	○	×
2019Boukaram	密	MV	H2	RSVD	×	×	○
2019Cambier	疎	LU	H2	RRQR	×	×	×
2019Zaspel	密	MV	H	ACA	×	×	○

# 並列化



R. Kriemann (2005), Parallel {H}-Matrix Arithmetics on Shared Memory Systems

R. Kriemann (2015), H-LU factorization on many-core systems



**JHPCN**

A. Ida  
I. Yamazaki  
S. Oshima  
T. Hiraishi

T. Iwashita  
K. Nakajima  
T. Aoki  
J. Dongarra

共有メモリ

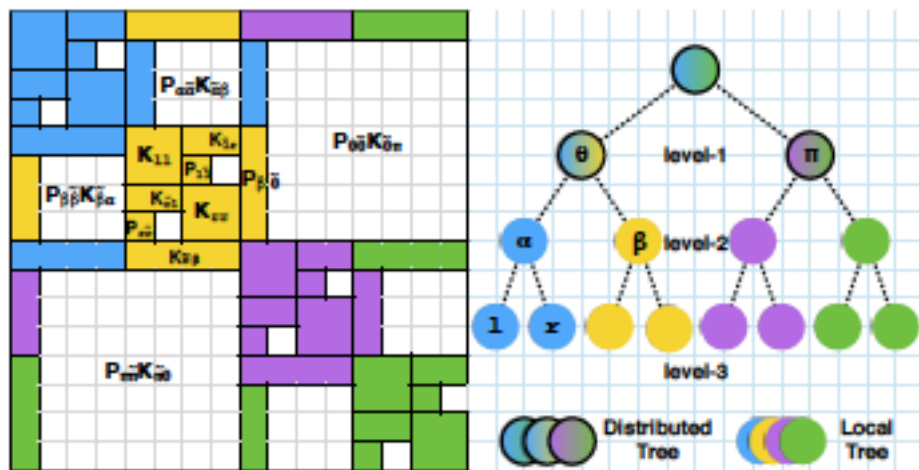
分散メモリ

M. Izadi (2012), Hierarchical Matrix Techniques on Massively Parallel Computers

S. Wang (2013), Efficient Scalable Algorithms for Solving Dense Linear Systems with HSS

Y. Li (2016), Distributed-memory Hierarchical Interpolative Factorization

C. D. Yu (2016), INV-ASKIT: A Parallel Fast Direct Solver for Kernel Matrices

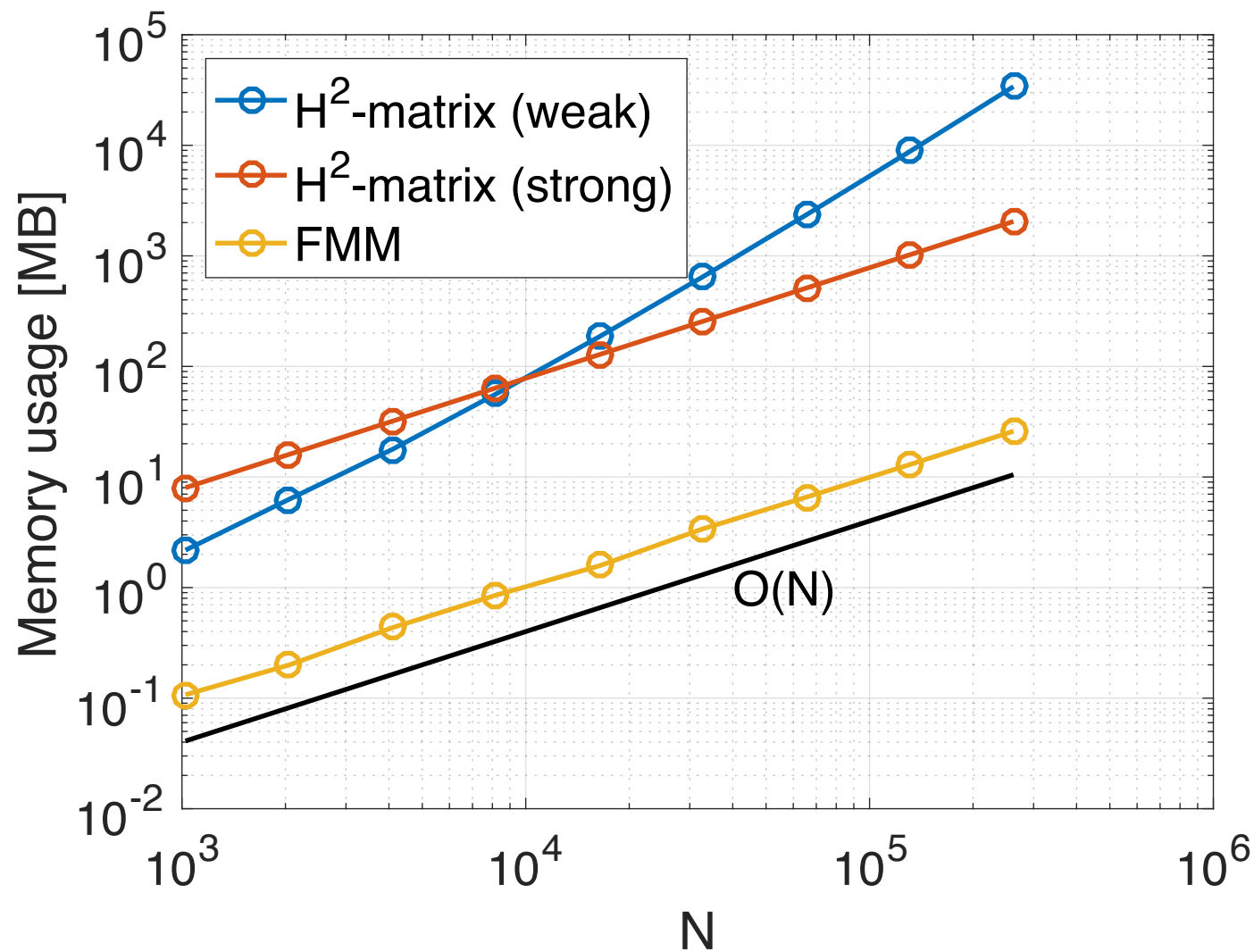


	Complexity	Concurrency
BLR	$\mathcal{O}(N^{4/3})$	High
H <sup>2</sup> (HSS)	$\mathcal{O}(N)$	Low

complexity-concurrency tradeoff



# FMMとH<sup>2</sup>-matrixの関係

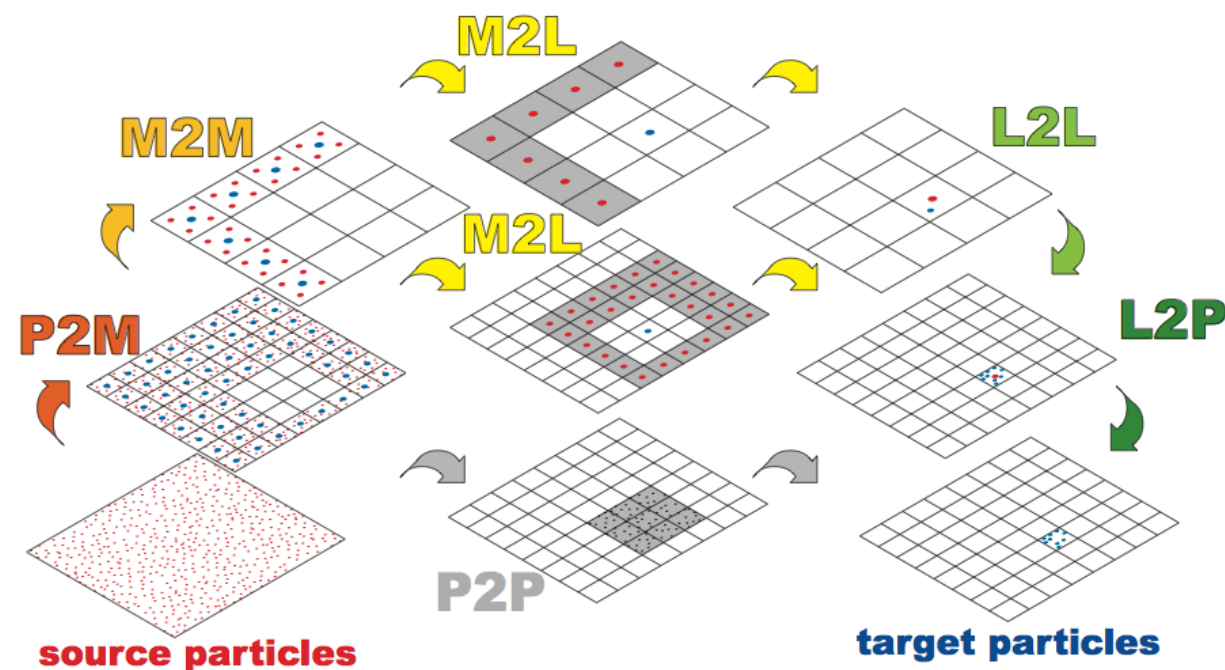
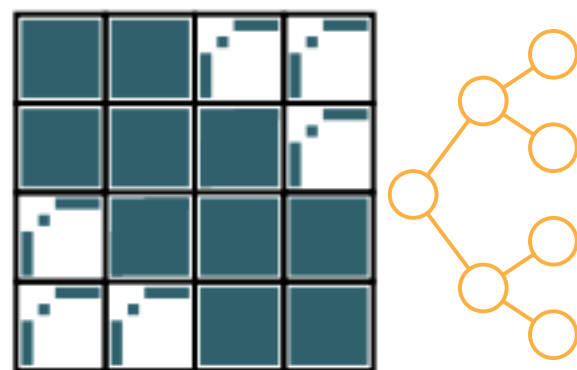


FMMはmatrix-freeのH<sup>2</sup>-matrix

FMMは動径基底関数しか扱えない

weak admissibility

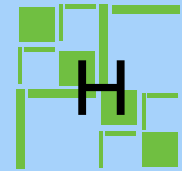
strong admissibility



# My new C++ code

C++ Class

dense low-rank hierarchical

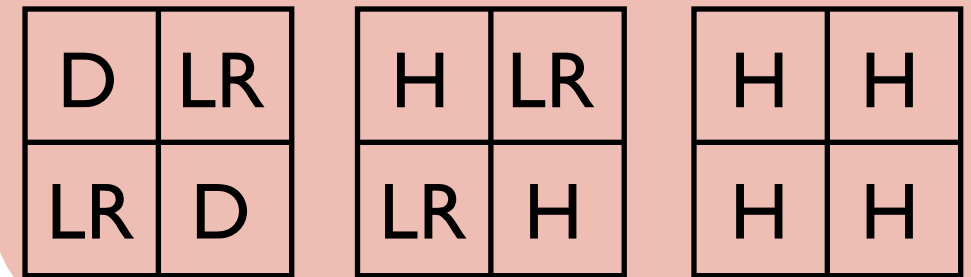


Operator overload

$D=D+D$   
 $L=D*L$   
 $H=H+H$

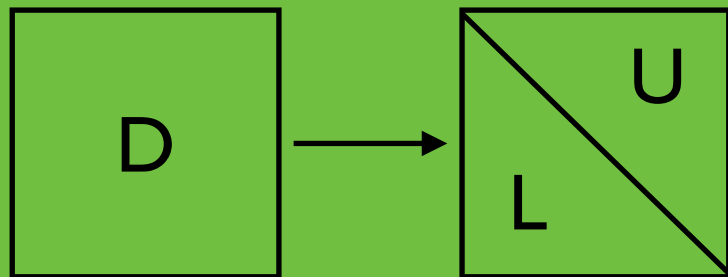
constructor  
 destructor

Dense  $D(N,N)$   
 LowRank  $LR(D,rank)$   
 Hierarchical  $H(2,2)$

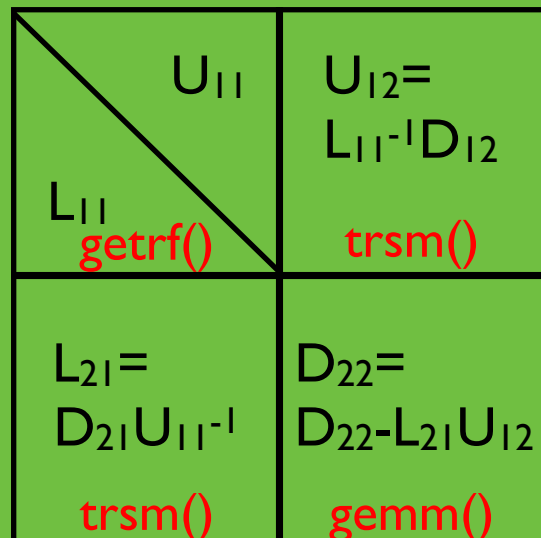


LU decomposition

$D.getrf()$



$H.getrf()$



$D.trsm(L)$

$$L_{11}^{-1}D_{12}$$

$H_D.trsm(H_L)$

$$\begin{bmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{bmatrix}^{-1} \begin{bmatrix} D_{11} & D_{12} \\ D_{21} & D_{22} \end{bmatrix}$$

$$= \begin{bmatrix} L_{11}^{-1}D_{11} & L_{22}^{-1}D_{12} \\ -L_{22}^{-1}L_{21}L_{11}^{-1}D_{11} + L_{22}^{-1}D_{21} & -L_{22}^{-1}L_{21}L_{11}^{-1}D_{12} + L_{22}^{-1}D_{22} \end{bmatrix}$$

$D.gemm(L,U)$

$$D_{22} - L_{21}U_{12}$$

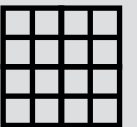
$H_D.gemm(H_L,H_U)$

$$\begin{bmatrix} D_{11} & D_{12} \\ D_{21} & D_{22} \end{bmatrix} = \begin{bmatrix} L_{11} & L_{12} \\ L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ U_{21} & U_{22} \end{bmatrix}$$

Any H structure

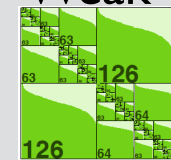
Hierarchical  $H(4,4)$

BLR too



admissibility

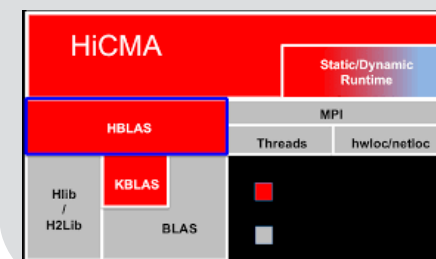
Weak



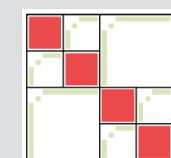
Strong



Hierarchical  $A(laplaceld, randx, N, N, rank, nleaf, admis, nblocks, nblocks);$



nested basis

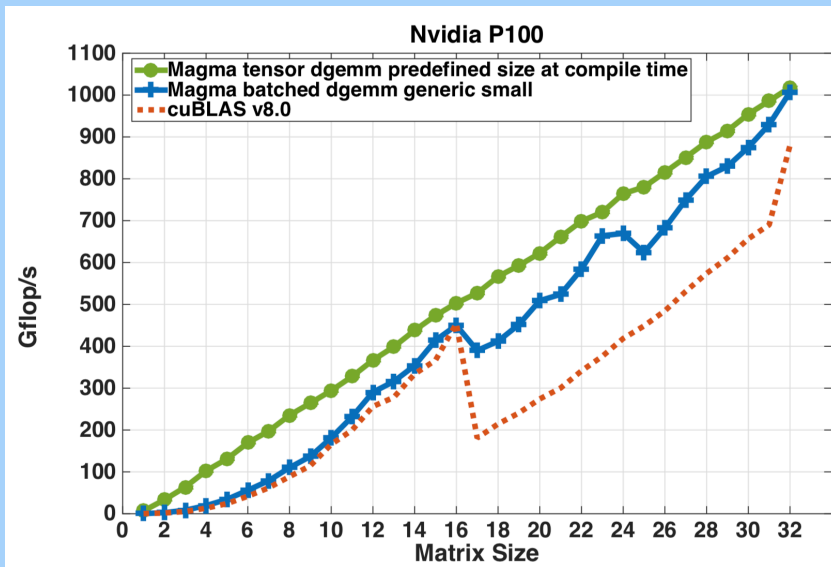
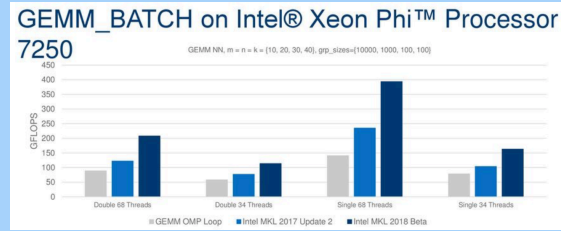
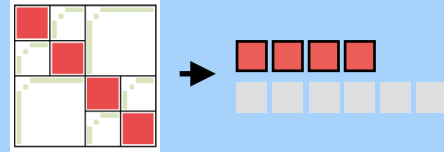




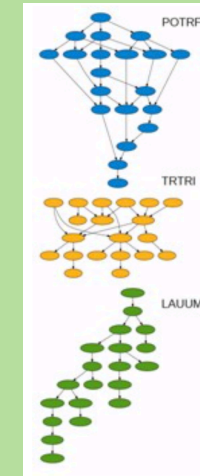
# GPU implementation

## Batch KBLAS

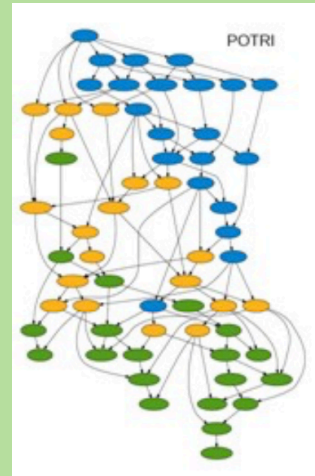
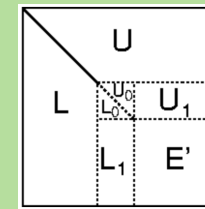
batch GEMM  
batch GEMV



## Runtime for LU

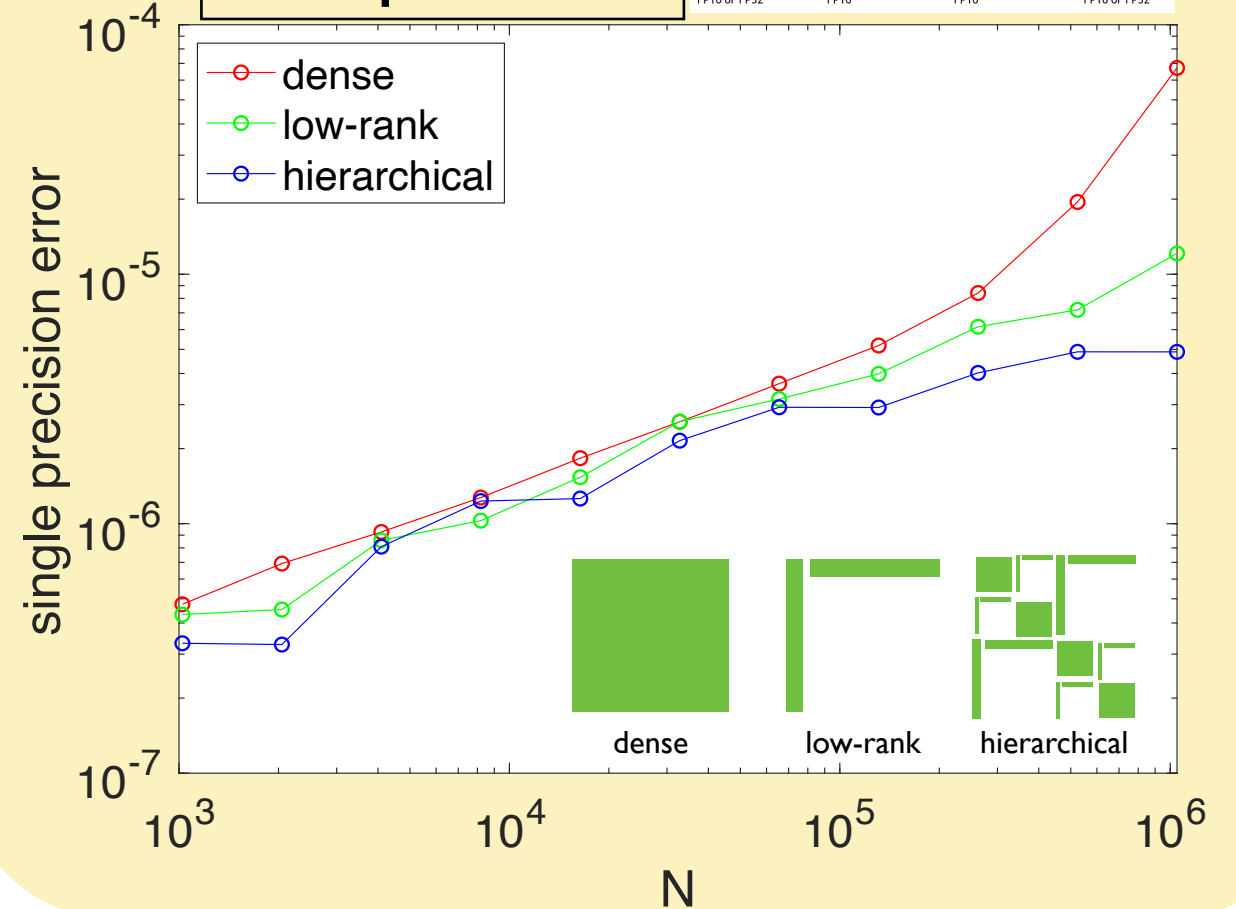


starPU  
OmpSs



## Low precision

$$D = \begin{matrix} \begin{matrix} A_{1,1} & A_{1,2} & A_{1,3} & A_{1,4} \\ A_{2,1} & A_{2,2} & A_{2,3} & A_{2,4} \\ A_{3,1} & A_{3,2} & A_{3,3} & A_{3,4} \\ A_{4,1} & A_{4,2} & A_{4,3} & A_{4,4} \end{matrix} & \begin{matrix} B_{1,1} & B_{1,2} & B_{1,3} & B_{1,4} \\ B_{2,1} & B_{2,2} & B_{2,3} & B_{2,4} \\ B_{3,1} & B_{3,2} & B_{3,3} & B_{3,4} \\ B_{4,1} & B_{4,2} & B_{4,3} & B_{4,4} \end{matrix} & + & \begin{matrix} C_{1,1} & C_{1,2} & C_{1,3} & C_{1,4} \\ C_{2,1} & C_{2,2} & C_{2,3} & C_{2,4} \\ C_{3,1} & C_{3,2} & C_{3,3} & C_{3,4} \\ C_{4,1} & C_{4,2} & C_{4,3} & C_{4,4} \end{matrix} \\ \text{FP16 or FP32} & \text{FP16} & & \text{FP16 or FP32} \end{matrix}$$



Contents lists available at [ScienceDirect](http://ScienceDirect)

**Parallel Computing**

journal homepage: [www.elsevier.com/locate/parco](http://www.elsevier.com/locate/parco)

Batched QR and SVD algorithms on GPUs with applications in hierarchical matrix compression

Wajih Halim Boukaram<sup>a,\*</sup>, George Turkiyyah<sup>b</sup>, Hatem Ltaief<sup>a</sup>, David E. Keyes<sup>a</sup>

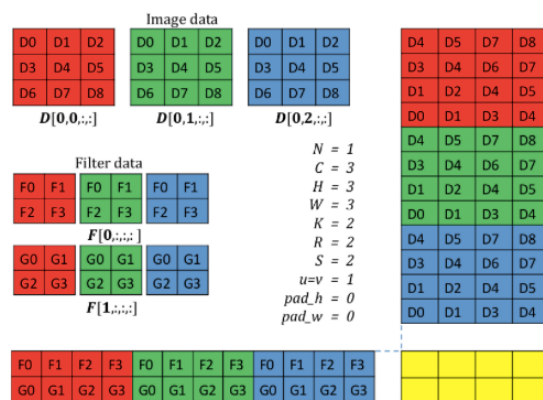
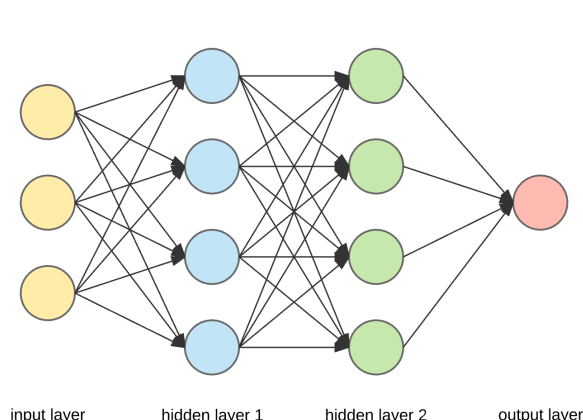
batch QR  
batch SVD  
batch RSVD  
batch ACA (variable M,N,K)

# H行列は深層学習に適用できるか？

深層ニューラルネットは  
密行列演算になる

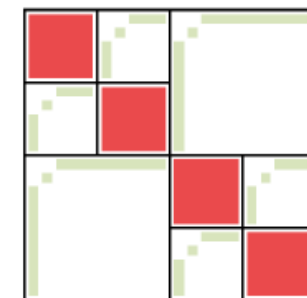
H行列は密行列の  
高速近似解法

畳み込みNNは  
密行列積になる

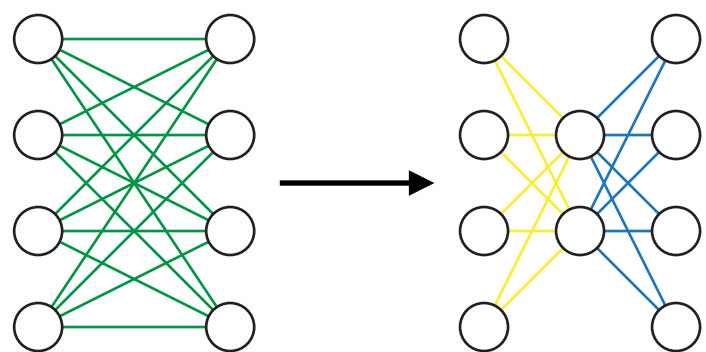


<https://arxiv.org/abs/1410.0759>

行列が非常に長細い  
ブロックのランクが高い  
圧縮した行列が1回しか使われない



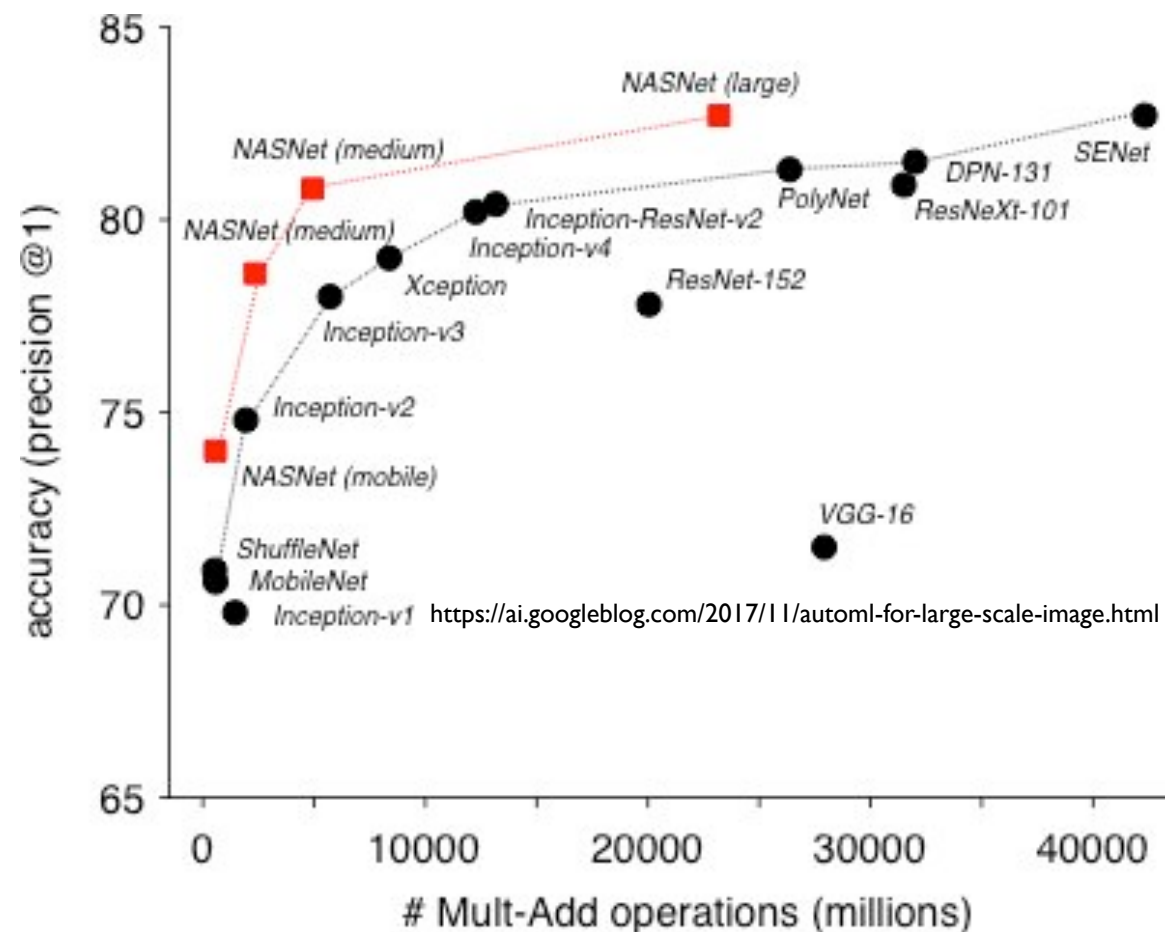
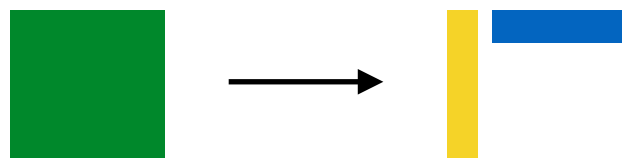
## 低ランク近似を用いてNN自体を圧縮



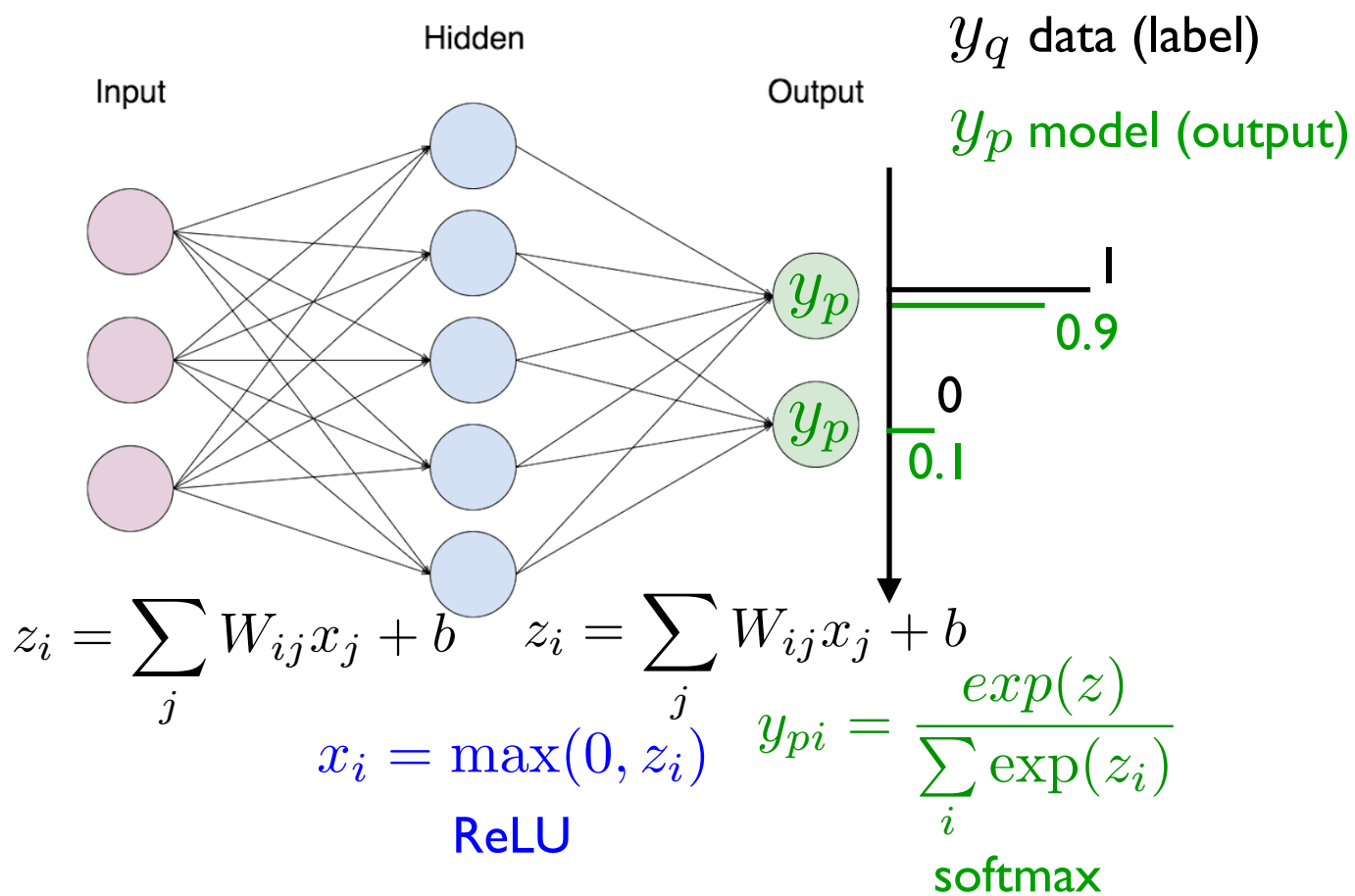
汎化性能と敵対的データに対する  
堅牢性が維持できるのか？

全結合層にしか有効でない

AutoMLの出現



# H行列とKronecker因子分解



Cross entropy loss function

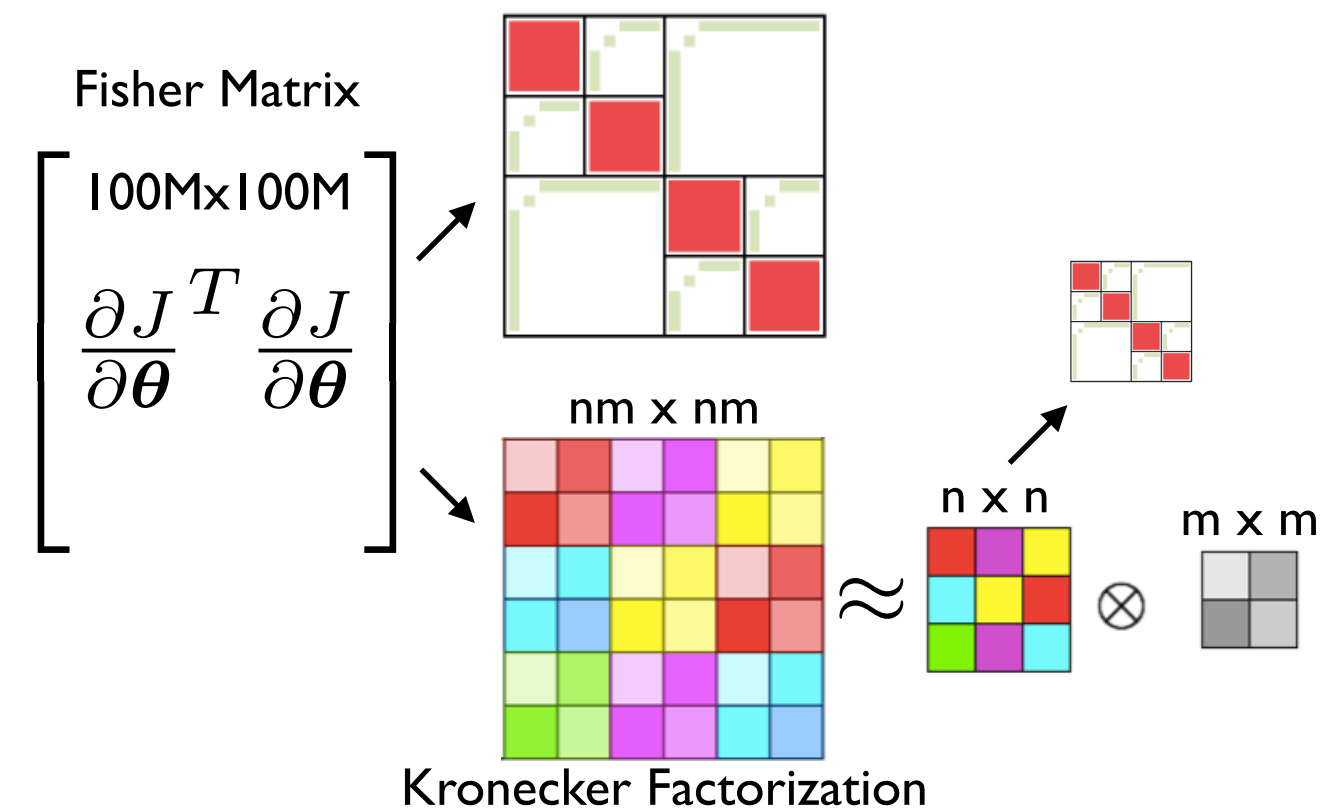
$$\begin{aligned}
 J(\theta) &= E_q(-\log p_\theta(y|\mathbf{x})) \\
 &= \sum_{(\mathbf{x}, y)} -q(y|\mathbf{x}) \log p_\theta(y|\mathbf{x}) \\
 &= \sum_{\mathbf{x}} \{-y_q \log y_p - (1 - y_q) \log(1 - y_p)\} \\
 &= \sum_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \theta)
 \end{aligned}$$

Back propagation

$$\begin{aligned}
 \frac{\partial J}{\partial W_{ij}} &= \frac{\partial J}{\partial y_{pi}} \frac{\partial y_{pi}}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} \\
 &= \sum_{\mathbf{x}} \frac{\partial \mathcal{L}}{\partial W_{ij}} = \sum_{\mathbf{x}} \frac{\partial \mathcal{L}}{\partial z_i} \frac{z_i}{W_{ij}} = \sum_{\mathbf{x}} g_i a_j
 \end{aligned}$$

Kronecker Product

Hierarchical Low-rank?



$$\begin{aligned}
 \left( \frac{\partial J^T}{\partial \theta} \frac{\partial J}{\partial \theta} \right)^{-1} &= \left\{ \left( \sum_{\mathbf{x}} \mathbf{g} \otimes \mathbf{a} \right)^T \left( \sum_{\mathbf{x}} \mathbf{g} \otimes \mathbf{a} \right) \right\}^{-1} \\
 &\approx \left( \sum_{\mathbf{x}} \mathbf{g}^T \mathbf{g} \right)^{-1} \otimes \left( \sum_{\mathbf{x}} \mathbf{a}^T \mathbf{a} \right)^{-1} \\
 &= G^{-1} \otimes A^{-1}
 \end{aligned}$$

# まとめ

- 行列の低ランク構造は元となる計算点の幾何学的な配置に依存
- 疎行列は計算点の接続, ランクは計算点の距離に関する
- FMM は matrix-free の $H^2$ 行列-ベクトル積
- 行列の形でFMMの演算を保存しておくことは左辺が多数ある場合に有効
- 階層的低ランク近似の手法間の違いは基底の共有と許容条件の違い
- Nullity theoremは $A$ と $A^{-1}$ の厳密なランク構造が同じことを保証するが数値的なランクは異なる
- 深層学習で扱う高次元の空間では低ランク近似よりもクロネッカー因子分解のほうが相性が良い