

# Hierarchical low-rank approximation methods on distributed memory and GPUs



## 背景

H行列、 $H^2$ 行列、HSS行列などの階層的低ランク近似法は $O(N^2)$ の要素を持つ密行列を $O(N)$ の要素を持つ行列に圧縮することができる。圧縮された行列を用いることで、行列積、LU分解、固有値計算を $O(N \log^2 N)$ で行うことができるため、従来密行列の解法が用いられてきた分野では階層的低ランク近似法の導入が近年盛んに行なわれている。また、密行列のみならず、疎行列の直接解法におけるSchur補元の圧縮に用いることもできるため、流体、構造、電磁界解析において前処理法として用いる研究も盛んに行なわれている。しかし、これらの階層的低ランク近似法は比較的新しい手法であるため、高性能な並列実装は少なく、GPUなどへの実装も未成熟である。これらの階層的低ランク近似法に内在する並列度は高く、高性能な分散メモリ・GPU実装に大きな期待が寄せられている。

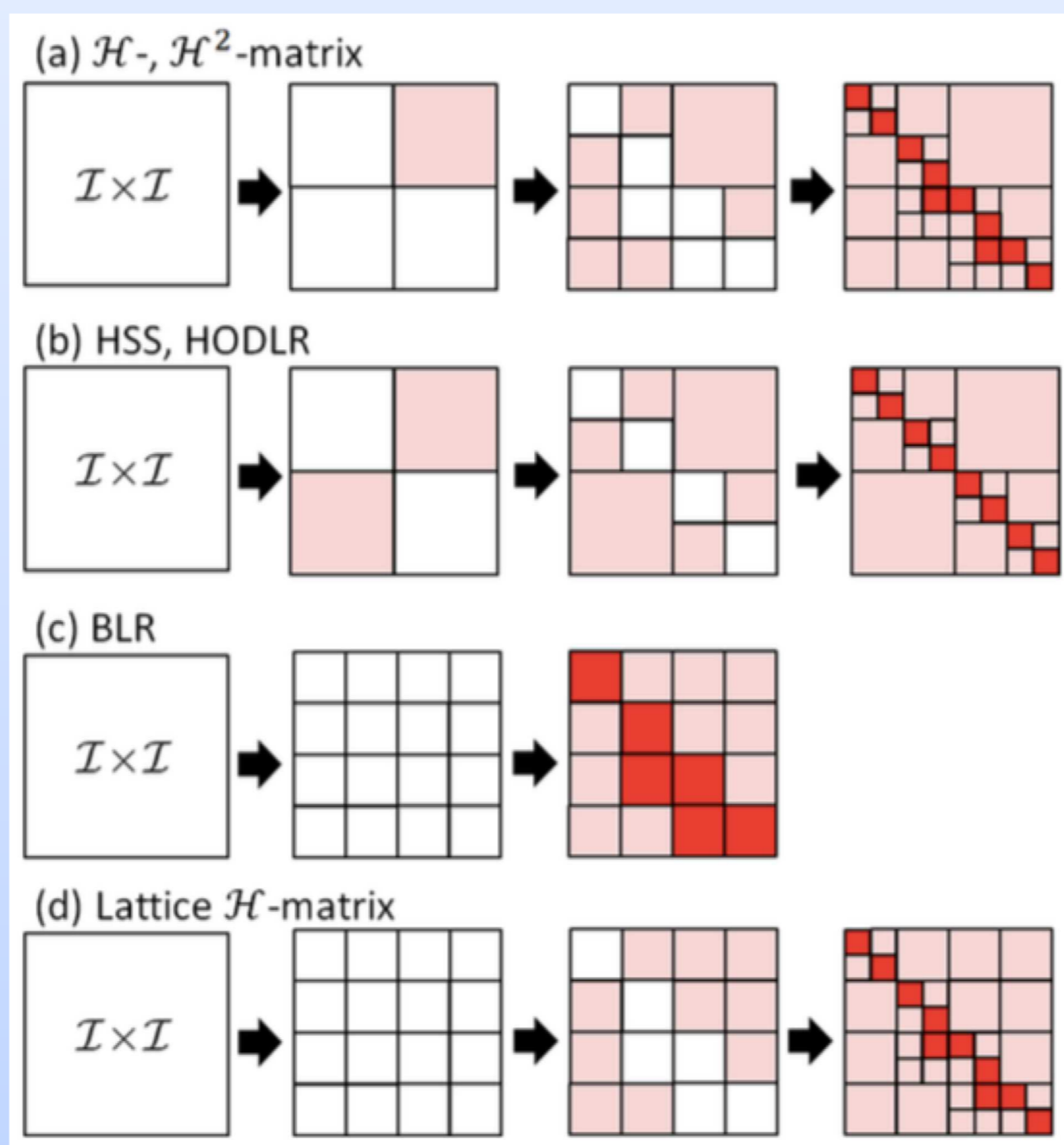
## 目的

本研究では、エクサスケールを視野に入れた階層的低ランク近似法の分散メモリ・GPU上での高性能な実装を行うことを目的とする。このとき重要になるのが比較的小さな密行列の高速な処理である。Tennessee大学のDongarraグループではまさにこのような小さな密行列のバッチ処理をGPU上で高速に行うライブラリを開発しており、JHPCNの国際共同研究として行うことでこの技術をいち早く導入できる。

昨年度は分散並列でH行列のLU分解を行う際の負荷分散アルゴリズムやバッチ処理を用いる際の階層的データのストリーム化について開発を行ったが、今年度はStarPUなどのランタイムを用いたH行列のLU分解やマルチGPU上でのH行列のLU分解、さらにH行列のQR分解も行う。

## Lattice H行列によるLU分解

右図に様々な低ランク近似法の違いを図示する。(a)にあるH-matrixや $H^2$ -matrixなどの手法は(b)のHSSやHODLRとは異なり、非対角ブロックをより細かく分割するのが特徴である。(c)のBLRは階層的でない低ランク近似法であり、(d)はBLRとH-matrixのハイブリッドである。(d)のlattice H-matrixはハイブリッド化

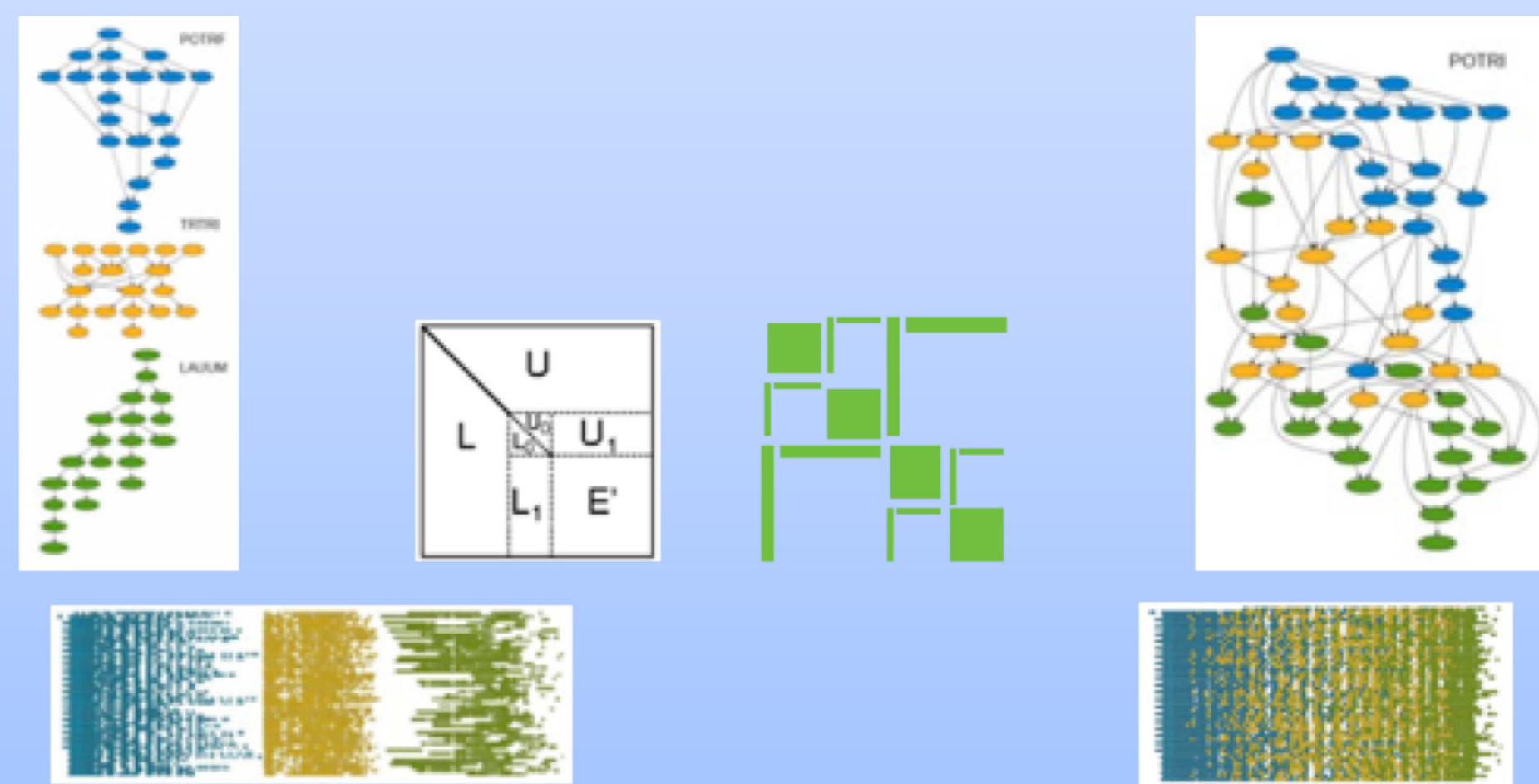


により、BLRのもつ並列度とH-matrixのもつ $O(N \log^2 N)$ の計算コストの両方を有する。下の表にLattice H行列の主要な関数の呼ばれる回数と計算コスト、その積である全体の計算コストを示す。

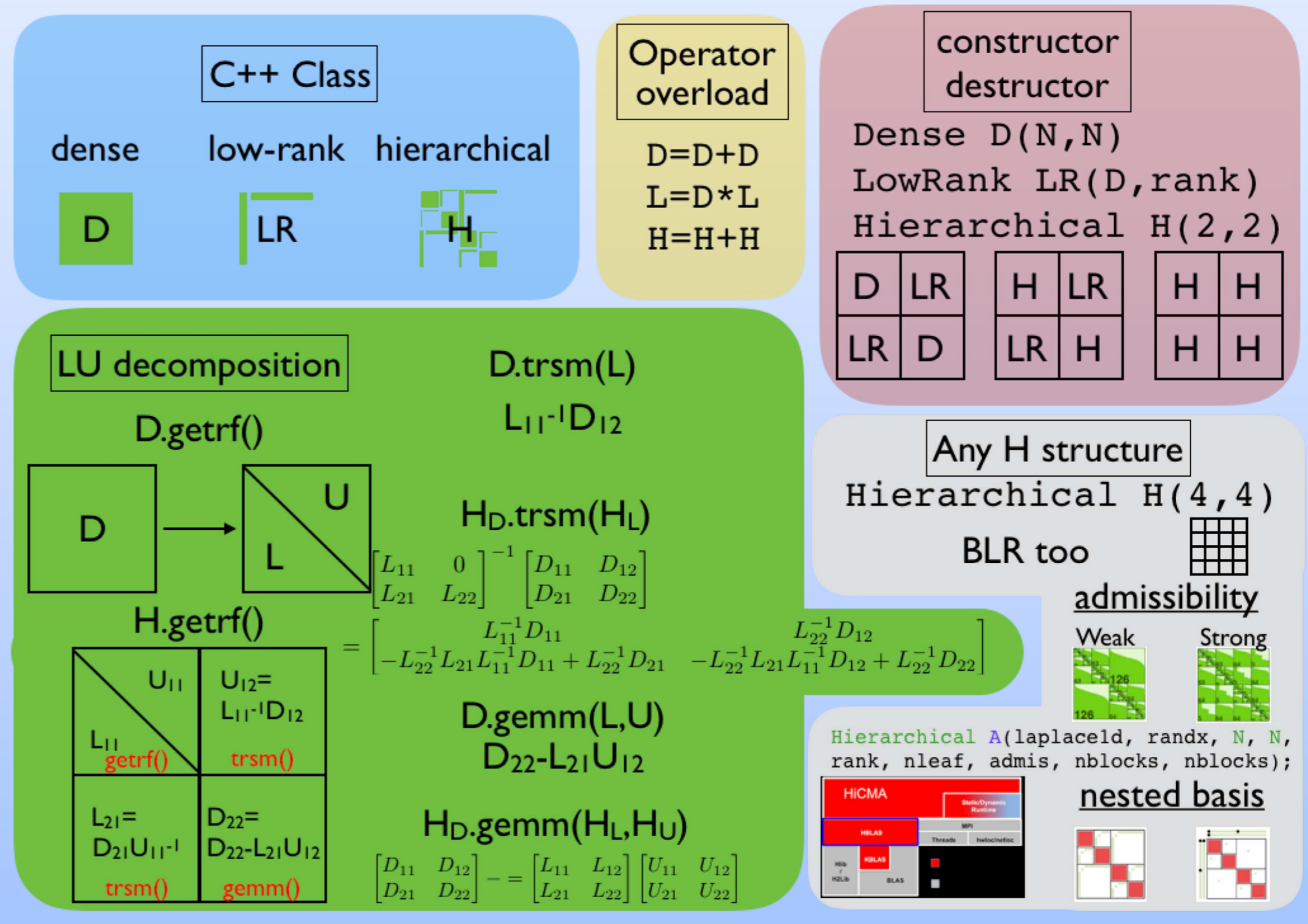
| Function                       | # of calls   | Complexity      | Total complexity            |
|--------------------------------|--------------|-----------------|-----------------------------|
| $\text{getrf}(H)$              | $O(n/l)$     | $O(l \log^2 l)$ | $O((n/l) * (l \log^2 l))$   |
| $\text{trsm}(\text{low-rank})$ | $O((n/l)^2)$ | $O(l)$          | $O((n/l)^2 * l)$            |
| $\text{trsm}(H)$               | $O(n/l)$     | $O(l \log l)$   | $O((n/l) * (l \log l))$     |
| $\text{gemm}(\text{low-rank})$ | $O((n/l)^3)$ | $O(l)$          | $O((n/l)^3 * l)$            |
| $\text{gemm}(H)$               | $O((n/l)^2)$ | $O(l \log^2 l)$ | $O((n/l)^2 * (l \log^2 l))$ |

## H行列のLU分解におけるランタイムの活用

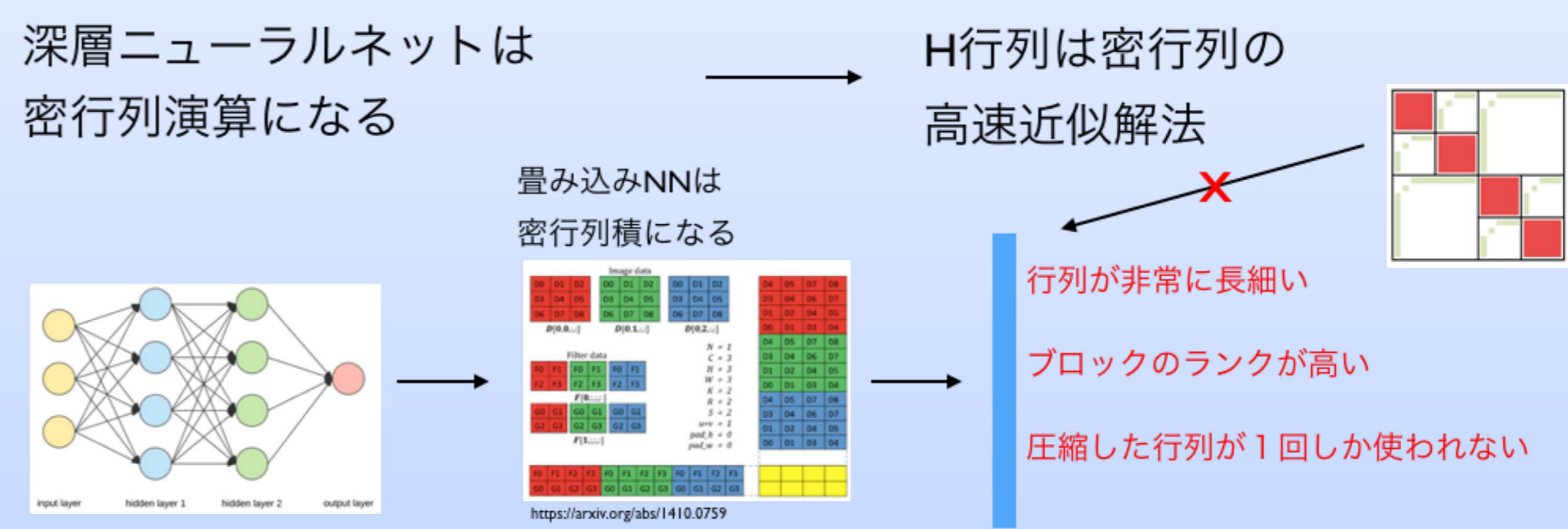
H行列を圧縮する計算やH行列・ベクトル積は個々のブロックを独立に計算できるため、依存関係のないバッチ処理になる。一方、LU分解では個々のブロックの計算が依存関係にあるため、単純なバッチ処理を行うことができない。密行列のブロック並列化では既に広く行われているランタイムによる依存関係の処理を利用した動的負荷分散は下図に示すように、ブロックLU分解などの処理において大幅な計算時間の短縮につながる。ただし、H行列の場合は、これに階層構造と、低ランクブロックという要素が追加されるため、動的負荷分散は容易ではない。



## 新たなC++コードの設計



## 深層学習におけるH行列の利用



## 低ランク近似を用いてNN自体を圧縮

