
固有値計算のための高性能 精度保証ライブラリの開発： 最新成果と自動チューニング機能

片桐 孝洋

名古屋大学情報基盤センター

第13回 自動チューニング技術の現状と応用に関するシンポジウム (ATTA2021)
2021年12月13日 (月) 15:15-15:45



名古屋大学
NAGOYA UNIVERSITY

2021年度 JHPCN採択課題(国際課題)

▶ タイトル

- ▶ Developing Accuracy Assured High Performance Numerical Libraries for Eigenproblems

▶ 期間

- ▶ 2021年4月～2022年3月 (1年間)
- ▶ 3年計画の3年目 (最終年度)

▶ 構成

▶ 副代表：

- ▶ Weichung Wang (National Taiwan U.),
Takeshi Ogita (Tokyo Women's U.)

▶ 主な構成員：

- ▶ K. Nakajima (U. of Tokyo), S. Ohshima (Nagoya U.),
K. Ozaki (Shibaura I.T.), O. Marques (LBNL),
▶ F-N. Hwang (National Central U.)



Aim of This Study

To establish accuracy assured numerical computing, we focus on the following three topics:

1. Developing **an accuracy assured numerical libraries for eigenproblems**;
2. Development of **high-performance implementation and auto-tuning technology for the developed accuracy assured numerical library**;
3. Discussing an **extension to non-linear problems** based on obtained knowledge of accuracy assured algorithms.



(1) Interdisciplinary Research

Researchers for Accuracy Assurance Problems

Researchers for High-performance Computing



The main target of 2021 FY (3) Developing Auto-tuning Function

(2) Assured algorithm

Ogita-Aishima's Iteration for Accurate Eigenpairs
(Quadratic Convergence)

A Result from 2020 FY

$$Ax = \lambda x, AX = XD$$

A: real symmetric matrix (binary64)

X: eigenvector matrix, D: eigenvalue matrix

Supercomputer "Flow" (Type I)

(matrix dimension: 2^{16} , 64 nodes)

pdsyevd: 240 sec, Iteration: 400 sec: **Ratio: 1.66**

Oakforest-PACS

(matrix dimension: 2^{16} , 64 nodes)

pdsyevd: 880 sec, Iteration: 2120 sec: **Ratio: 2.40**

Accurate matrix multiplication is performed by Ozaki Method with six matrix multiplications.

- 1: $R \leftarrow I - \widehat{X}^T \widehat{X}$
- 2: $S \leftarrow \widehat{X}^T A \widehat{X}$
- 3: $\tilde{\lambda}_i \leftarrow \frac{s_{ii}}{1 - r_{ii}}$ (for $1 \leq i \leq l$)
- 4: $\tilde{D} \leftarrow \text{diag}(\tilde{\lambda})$
- 5: $\omega \leftarrow 2(\|S - \tilde{D}\|_2 + \|A\|_2 \|R\|_2)$
- 6: $\tilde{e}_{ij} \leftarrow \begin{cases} \frac{s_{ij} + \tilde{\lambda}_j r_{ij}}{\tilde{\lambda}_j - \tilde{\lambda}_i} & (|\tilde{\lambda}_i - \tilde{\lambda}_j| > \omega) \\ r_{ij}/2 & (\text{otherwise}) \end{cases}$ (for $1 \leq i, j \leq l$)
- 7: $\tilde{X} \leftarrow \widehat{X} + \widehat{X} \tilde{E}$

Lower is better.

Adaptation of Supercomputers



(Ref. 内野 佑基, 尾崎 克久, 荻田 武史: "実対称行列の固有値分解に対する反復改良法の大規模並列環境における実装と評価", 第177回ハイパフォーマンスコンピューティング研究発表会, Dec. 2020)

VNC (Verified Numerical Computation)-HPC

ライブラリ

- **ポスト「京」萌芽的課題1：基礎科学のフロンティアー極限への挑戦ー：「極限の探究に資する精度保証付き数値計算学の展開と超高性能計算環境の創成」**（代表：荻田武史 東京女子大学教授）～2019年（終了）で開発
- 精度が保証された計算結果を実用的に得られるような超高性能計算環境を構築

- 以下のライブラリのソースコードを公開中

1. **OzBLAS: Accurate and Reproducible BLAS based on Ozaki scheme** (尾崎スキームに基づく高精度基本線形計算ライブラリ) [for PC, GPU]
2. **GEMMTC: GEMM using Tensor Cores** (Tensorコアを用いた行列積プログラム) [for GPU]
3. **DHPMM_F for GPU: High-precision Matrix Multiplication with Faithful Rounding** (高精度行列積プログラム) [for GPU]
4. **PDDOTK: K-fold Precision Dot Product** (分散並列版高精度内積計算プログラム) [for PC, FX100]
5. **BLAS-DOT2: Higher-precision BLAS based on Dot2** (高精度内積計算アルゴリズムDot2に基づく基本線形計算ライブラリ) [for GPU]
6. **LINSYS_VR: Verified Solution of Linear Systems with Directed Rounding** (連立一次方程式に対する精度保証プログラム(丸めモードの変更機能付き)) [for K Computer, FX100]
7. **LINSYS_V: Verified Solution of Linear Systems** (連立一次方程式に対する精度保証プログラム) [for PC, K Computer, FX100]
8. **DHPMM_F: High-precision Matrix Multiplication with Faithful Rounding** (高精度行列積プログラム) [for PC, K Computer, FX100]

HP : <http://www.math.twcu.ac.jp/ogita/post-k/index.html>

Verified Numerical Computations
VNC - HPC
High-Performance Computing

ポスト「京」萌芽的課題1
基礎科学のフロンティアー極限への挑戦
極限の探究に資する精度保証付き数値計算学の展開と超高性能計算環境の創成

トップページ Top	課題概要 Abstract	研究体制 System	研究内容 Contents	成果公開 Results
---------------	------------------	----------------	------------------	-----------------

リンク

- トップページ
- 課題概要
- 研究体制
- 研究内容
- 成果公開

課題概要

本研究課題の目的は、ポスト「京」において、精度が保証された計算結果を実用的に得られるような超高性能計算環境を構築することです。具体的には、ポスト「京」で実行される様々な数値シミュレーションにおいて、数値計算による計算誤差の問題が解消されることによって、シミュレーションサイエンスの品質を向上させ、さらに想定外の現象が発生する可能性を低減することが可能となります。すなわち、本研究課題の遂行は、人が安心して生活できる社会基盤の構築に直結します。

たとえば、産業界における製品開発や非破壊検査等のための計算工学シミュレーションや、地震や津波などの災害シミュレーションは、日本に限らず世界中で盛んに行われていますが、計算誤差の観点から高精度なシミュレーションの実現をしている例は皆無です。これは、計算機によって問題の近似的な解を得ることよりも、その近似解の検証のほうがはるかに困難かつ計算資源を必要とする、と考えられており、それが高性能計算分野の常識であるからです。実際、これは1990年代までは事実でしたが、本研究グループが2000年代から切り拓いてきた高速で実用的な精度保証付き数値計算法やエラーフリー変換に基づく高精度数値計算法をベースとして、今、ポスト「京」によって、この常識を打ち破る時期が到来しています。すなわち、ポスト「京」において、高速性と高精度性を融合した超高性能計算環境の創成を世界に先駆けて達成することは、スーパーコンピュータに質的転換をもたらし、我が国の高い科学技術力を国内外に示すこととなります。

性能 = 速度 × 精度

大規模数値計算を行う高性能計算分野においては、問題の大規模化に伴って計算誤差が累積しやすくなり、これが今後、大きな問題となってきます。たとえば、計算機上で標準的に用いられる32ビットの浮動小数点演算では、100万次元程度の密行列系線形問題に対して数値計算を行うと、問題自身は比較的良条件下で解きやすい問題であったとしても、誤差解析の結果から理論的には1桁も正しくないような計算結果が得られることが分かっています。また、計算誤差の単なる累積だけでなく、問題の困難さが条件数として

名古屋大学 大学院情報学研究科 M1 青木 将太
名古屋大学情報基盤センター 准教授 大島聡史

精度保証ライブラリの自動チューニング
(FY2021の成果)



発表の流れ

- ▶ 背景
- ▶ 機械学習モデルの特徴量解析
- ▶ 高精度行列-行列積アルゴリズム
(尾崎の方法)
- ▶ 予備評価
- ▶ おわりに

発表の流れ

- ▶ **背景**
- ▶ 機械学習モデルの特徴量解析
- ▶ 高精度行列-行列積アルゴリズム
(尾崎の方法)
- ▶ 予備評価
- ▶ おわりに

背景

- ▶ 人工知能(AI)が出力する答えを検証が不十分なまま利用することで様々な社会的な問題が生じている
- ▶ AI出力に対して人間による検証が必須
- ▶ 性能チューニング工数の削減を目的に、ソフトウェア自動チューニング (Software Auto-tuning, AT) 技術が研究されている
 - ▶ AIのAT機構への適用が進んでいる
 - 精度保証数値計算ライブラリの性能パラメータチューニングにAIを適用した事例の<AI出力結果の説明性> を検証

→ Scientific XAI (SXAI)

第13回 自動チューニング技術の現状と応用に関するシンポジウム

発表の流れ

- ▶ 背景
- ▶ 機械学習モデルの特徴量解析
- ▶ 高精度行列-行列積アルゴリズム
(尾崎の方法)
- ▶ 予備評価
- ▶ おわりに

説明可能なAI (XAI)

- ▶ 機械学習が出した結果が説明できるか？
 - ▶ Explainable AI (XAI)
- ▶ 説明可能なAIには以下の2種がある[1]
 - ▶ 説明可能性
 - ▶ 予測の判断理由を人間が理解できるように説明する技術
 - ▶ 例) LIME、SHAPなどのツール
 - ▶ 解釈可能性
 - ▶ 内部構造を解析することで、予測に至る計算過程を確認できるようなAI技術
 - ▶ 例) 決定木を作る

→ここでは<説明可能性>を扱う

[1] 大坪ほか：「XAI(説明可能なAI)：そのとき人工知能はどう考えたのか」、リックテレコム、2021

大局説明と局所説明

▶ 大局説明 (Global Explanations)

- ▶ 「AIモデルの全体的なふるまい」を理解する[1]
- ▶ 例) SHAPなどのツール

▶ 局所説明 (Local Explanations)

- ▶ 「個々の予測結果の判断理由」を理解する[1]
- ▶ 例) LIMEなどのツール

→ ここでは、数値計算処理のATに適用したAI
について、双方の説明性を取り扱う

[1] 大坪ほか：「XAI(説明可能なAI)：そのとき人工知能はどう考えたのか」、リックテレコム、2021

LIME (Local Interpretable Model-agnostic explanations) [2]

- ▶ Local surrogate model
 - ▶ ブラックボックスモデルの個々の予測を説明するために用いられる解釈可能なモデル
- ▶ 理由を人間でも理解できるように提示
- ▶ 各特徴がどの程度分類に貢献しているか調べ分類器の予測結果を説明
 - ▶ 説明対象データを摂動させた近傍データからAIモデルを得る
- ▶ 任意の分類器に適用できる
- ▶ 欠点
 1. 毎回、説明がぶれる (乱数利用)
 2. ハイパーパラメタ調整がある
 3. 説明できないケースがある (→あきらめて他の方法を利用)

[2] M. T. Ribeiro, S. Singh, and C. Guestrin: Why should I trust you?: Explaining the predictions of any classifier, Proc. of 22nd ACM SIGKDD, pp.1135-1144, 2016.

SHAP (SHapley Additive exPlanations) [3]

- ▶ 協力ゲーム理論のシャープレイ値 (Shapley Value) を機械学習に応用
 - ▶ それなりの妥当性がある
- ▶ 近似的にシャープレイ値を算出
 - ▶ ツリー系アンサンブルモデル：高速で正確なシャープレイ値の算出
 - ▶ ディープラーニングモデル：高速で近似的なシャープレイ値の算出
 - ▶ 一般的なアルゴリズム：シャープレイ値の推定値の算出
- ▶ 欠点
 1. 計算量が多い (→近似的な評価値を利用)

[2] S. Lundberg, S-I. Lee, A Unified Approach to Interpreting Model Predictions, 2017
<https://arxiv.org/abs/1705.07874>

第13回 自動チューニング技術の現状と応用に関するシンポジウム



発表の流れ

- ▶ 背景
- ▶ 機械学習モデルの特徴量解析
- ▶ 高精度行列-行列積アルゴリズム
(尾崎の方法)
- ▶ 予備評価
- ▶ おわりに

Background, objective, and motivation of our work

- **High-accuracy and low-accuracy calculations** are one of the important calculation techniques which can solve large-scale and complicated calculations.
- Some studies have investigated the **accuracy assurance of BLAS and LAPACK**.
 - These studies have used mixed procedure computations and arbitrary digit computations.
 - Most of BLAS and LAPACK libraries tend to place **less emphasis on the accuracy of the computation results**.
- Why high-accuracy and assured calculations are not widespread?
- => **TIME**
- We consider to **calculate it on GPU and shorten the time**.
- This work will be helpful for large-scale and complicated calculations on current and future generation computers. (⇒ **topics of interest** of this workshop)

In particular, we will focus on the following **topics of interest**, but not limited to:

- Programming models, languages and frameworks for facilitating HPC software evolution and refactoring.
- Algorithms and implementation methodologies for future-generation computing systems, including manycores and accelerators (GPUs, Xeon Phi, etc).
- Automatic performance tuning techniques, runtime systems and domain-specific languages for hiding the complexity of underlying system architectures.
- Practices and experiences on porting of legacy applications and libraries.

High-precision matrix-matrix multiplication algorithm

- Our target calculation: assured matrix-matrix multiplication (MMM) method proposed by Ozaki et al.
 - hereinafter, we refer to this MMM method as **the Ozaki method**
- Overview of the Ozaki method:

consider $C = AB$

- A : a matrix of size $m * l$
- B : a matrix of size $l * n$
- C : a matrix of size $m * n$

Step1: error-free transformation

$$A = A^{(1)} + A^{(2)} + A^{(3)} + \dots + A^{(p)}$$

$$B = B^{(1)} + B^{(2)} + B^{(3)} + \dots + B^{(q)}$$

The elements in the matrices with lower indices are given with a higher number of digits.

Step2: individual MMM

$$\begin{aligned} AB &= (A^{(1)} + A^{(2)} + \dots + A^{(p)})(B^{(1)} + B^{(2)} + \dots + B^{(q)}) \\ &= \underline{A^{(1)} B^{(1)}} + \underline{A^{(1)} B^{(2)}} + \underline{A^{(2)} B^{(1)}} + \dots + \underline{A^{(p)} B^{(q)}} \end{aligned}$$

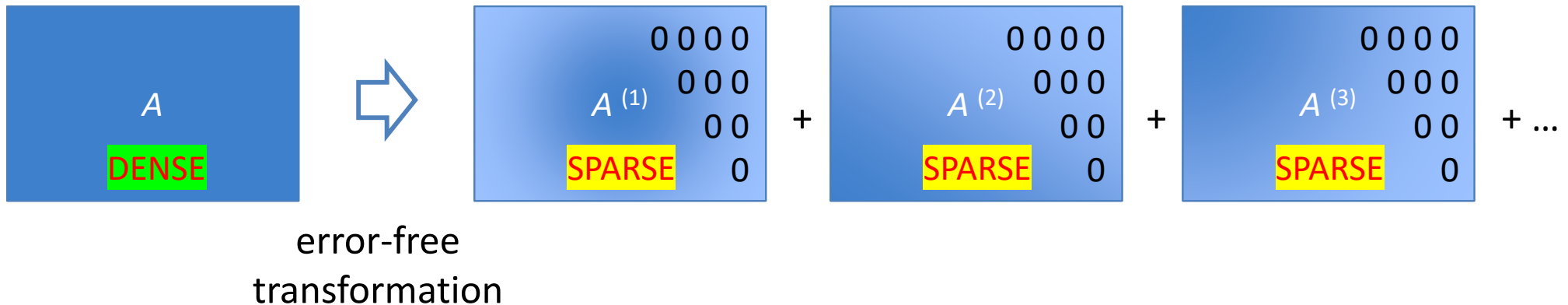
Step3: Accurate Sum

$$\begin{aligned} fl(A^{(i)} B^{(j)}) &= A^{(i)} B^{(j)} \text{ for } 1 \leq i \leq p, 1 \leq j \leq q. \\ &= fl(A^{(1)} B^{(1)}) + fl(A^{(1)} B^{(2)}) + fl(A^{(2)} B^{(1)}) + \dots \\ &\quad + fl(A^{(p)} B^{(q)}) \\ &= C_1 + C_2 + \dots + C_{pq} \end{aligned}$$

fl is a floating-point arithmetic with rounding to the nearest

Previous work and Proposed method

- When the value ranges of the input matrix elements are large, error-free transformation generates many sparse matrices.



- In this case, many double precision general matrix - matrix multiplication (dgemm) are performed. Transforming dense matrices into sparse matrices and **performing sparse matrix operations will require shorter calculation time than dense matrix operations.**
- Therefore, our previous work proposed to transform the target matrices into sparse matrices and calculate sparse matrix computations **on CPU.**
 - Considering the performance, sparse matrix - vector multiplications (SpMV) are used.
- In this study, we propose to calculate these sparse matrix computations **on GPU.**

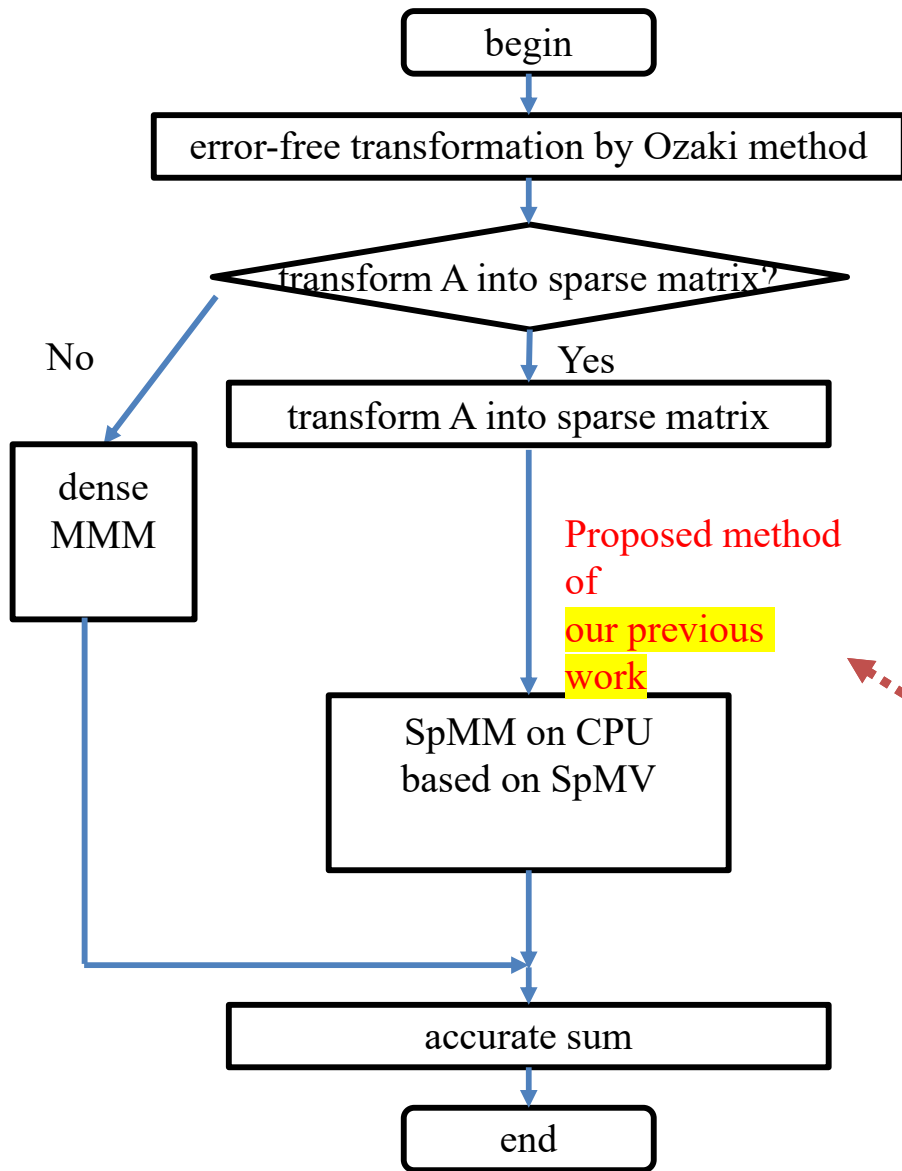
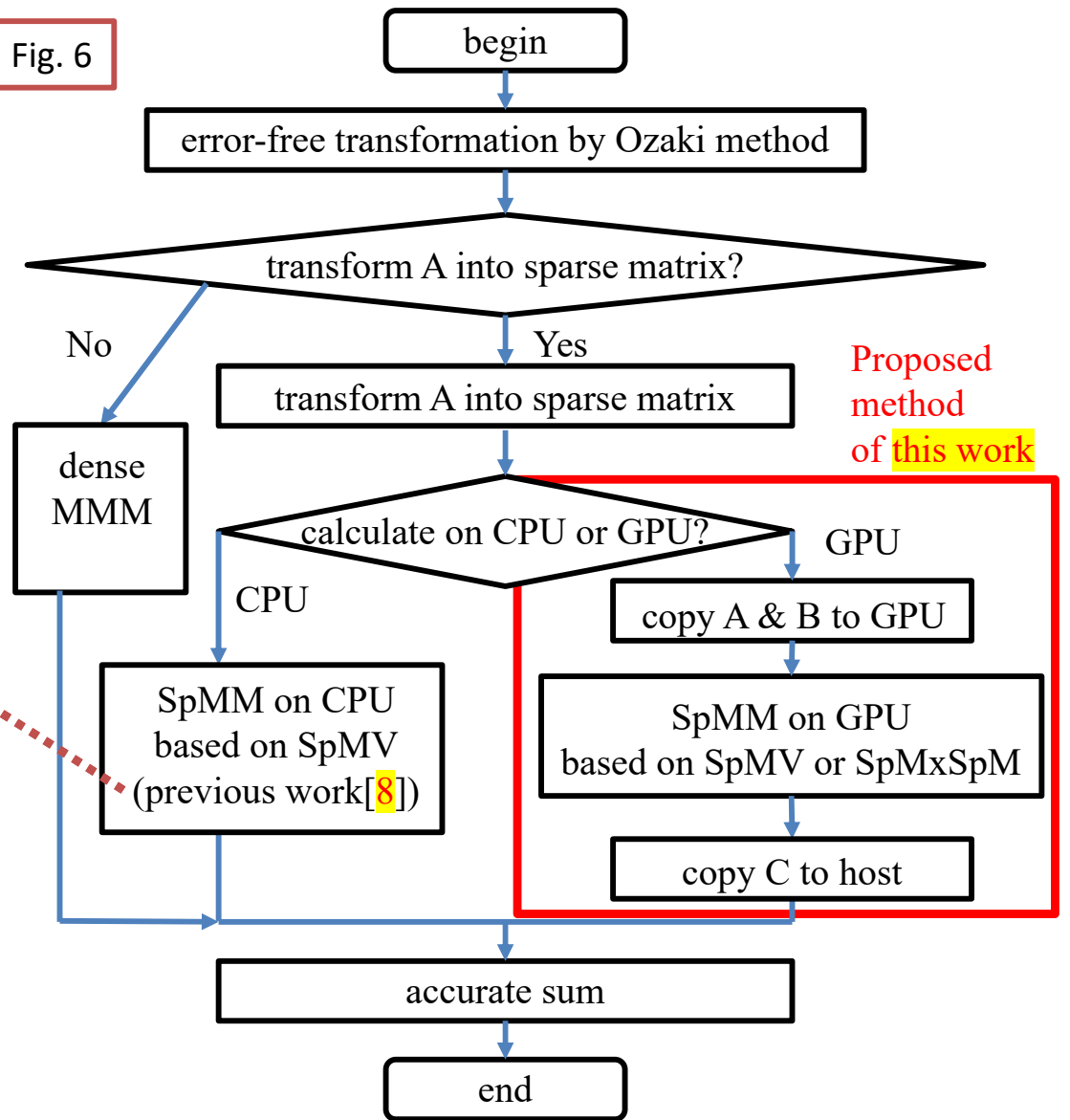
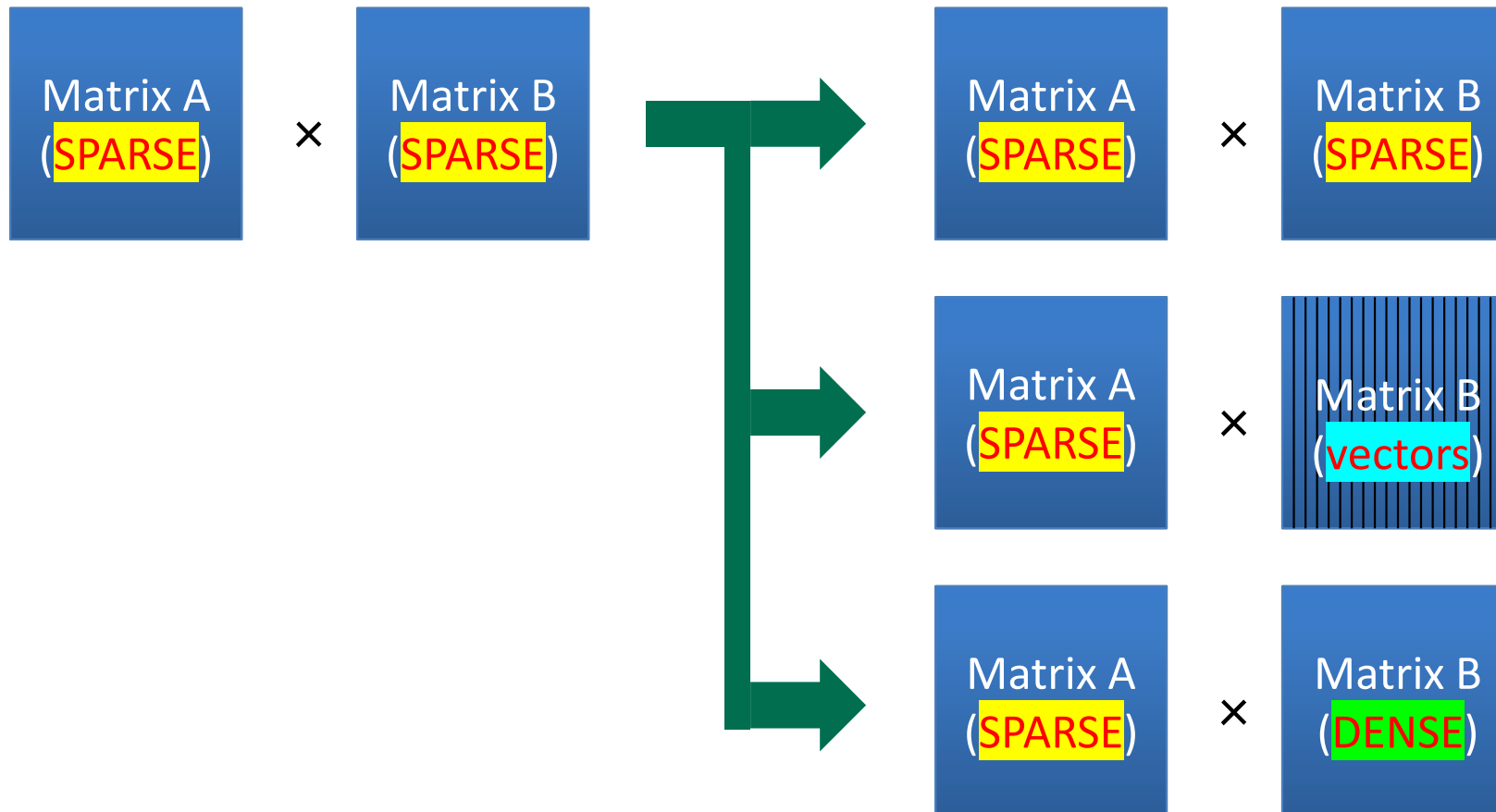


Fig. 6



How to calculate Sparse Matrix A - Matrix B multiplication?



高精度行列積ライブラリ(尾崎の方法)の実装種別と表記

1. dgemmによる実装 1dgemm
2. CRS形式でのSpMV(内部並列) 2CRS
3. CRS形式でのSpMV(外部並列) 3CRS
4. CRS形式でのSpMV(複数右辺による内部並列化) 4CRS
5. CRS形式でのSpMV(複数右辺による内部並列化(ブロッキング)) 5CRS
6. ELL形式でのSpMV(内部並列) 6ELL
7. ELL形式でのSpMV(外部並列) 7ELL
8. ELL形式でのSpMV(複数右辺による内部並列化) 8ELL
9. ELL形式でのSpMV(複数右辺による内部並列化(ブロッキング)) 9ELL
10. Batched BLASを適用した方式 10GPU
11. dgemmによる実装(GPU) 11GPU
12. CRS形式でのSpMV(GPU) 12GPU
13. ELL形式でのSpMV(GPU) 13GPU
14. CRS形式でのSpMM(GPU) 14GPU

発表の流れ

- ▶ 背景
- ▶ 機械学習モデルの特徴量解析
- ▶ 高精度行列-行列積アルゴリズム
(尾崎の方法)
- ▶ 予備評価
- ▶ おわりに

評価環境



▶ 計算機：名古屋大学情報基盤センター設置の スーパーコンピュータ「不老」Type II サブシステム

- ▶ CPU
 - ▶ Intel Xeon Gold 6230, 20コア, 2.10 - 3.90 GHz × 2 ソケット
- ▶ GPU：機械学習モデル生成用
 - ▶ NVIDIA Tesla V100 (Volta) SXM2, 2,560 FP64コア, upto 1,530 MHz × 4 ソケット
- ▶ LIMEは ver. 0.2.0.1 を利用
- ▶ SHAPは ver.0.39.0 を利用
- ▶ 機械学習モデル（分類器）
 - ▶ scikit-learn ver. 0.24.1のランダムフォレスト
 - ▶ 尾崎の方法の実装選択中の11種類

テスト行列生成方法

- ▶ **行列 1** : 行列 A, B の要素を 0~1 の範囲で生成し, ある疎度分の要素に対し $\text{pow}(10, \text{rand}() \% \Phi)$ で生成した値を挿入
 - ▶ $\Phi=30$ まで
(Φ が一定数を超えるとpythonで扱える範囲を超える)
- ▶ **行列 2** : 行列 A, B を単位行列とある疎度分の要素に対して 0~1 の範囲で生成
 - ▶ 疎度90から98
- ▶ **行列 3** : 行列 A, B を単位行列とある疎度分の要素に対して $\text{pow}(10, \text{rand}() \% \Phi)$ で生成した値を挿入
 - ▶ $\Phi=30$ まで
- ▶ **乱数行列 : 以上の乱数の種は固定**

学習データの数と予測精度

▶ 学習データ

1. 行列1: 大きな値が作れる、疎度がほぼ0
 2. 行列2: 0~1の範囲の疎行列
 3. 行列3: 大きな値が作れる、任意の疎度
- ▶ 行列サイズ1000~4000

▶ 機械学習情報

▶ 方法: ランダムフォレスト

▶ 説明変数: 7変数

- ▶ 1) 行列サイズ、2) 入力行列の疎度、3) Aの要素最大値、4) Aの要素最小値、5) A疎な分解数、6) A密な分解数、7) Bの分解数

▶ 学習データ数: 199個

▶ テストデータ数: 23個

▶ 正答率: 91.3%

※ここでは、汎用的で予測精度の高い分類器を作るのが目的ではないので、精度にはこだわらない

LIMEによる出力例

●あるケースによる予測結果を説明(局所説明)

Prediction probabilities

4CRS	0.99
1dgemm	0.01
3CRS	0.00
5CRS	0.00
Other	0.00

NOT 4CRS

4CRS



Feature Value

sp percentage	98.00
A max	1.00
sparse A	3.00

[1.5e+03 9.8e+01 1.0e+00 0.0e+00 3.0e+00 0.0e+00 4.0e+00]

- 予想される目的変数の確率を表示
「4CRS」が99%予測される

- LIMEによる説明変数の寄与度を表示
「4CRS」を予測する場合、「sp percentage」の寄与(正)が34%

- このケースの、説明変数の値を表示

LIMEによる分析結果(1/3)

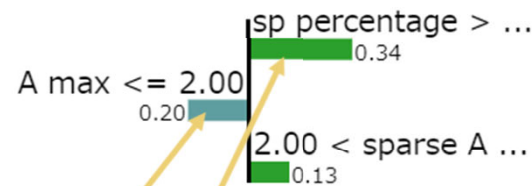
4CRSを99%予測

Prediction probabilities

4CRS	0.99
1dgemm	0.01
3CRS	0.00
5CRS	0.00
Other	0.00

NOT 4CRS

4CRS



Feature Value

sp percentage	98.00
A max	1.00
sparse A	3.00

[1.5e+03 9.8e+01 1.0e+00 0.0e+00 3.0e+00 0.0e+00 4.0e+00]

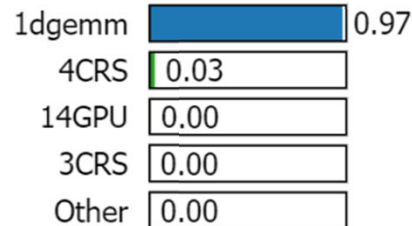
- 入力行列の疎度が高い(sp percentage)と4CRSと予想される
→ 取得データ・数値アルゴリズム的に、正しい判断
- Aの要素の最大値が小さいと別の実装が速い可能性がある
→ A maxが負の要因となってもおかしくない



LIMEによる分析結果(2/3)

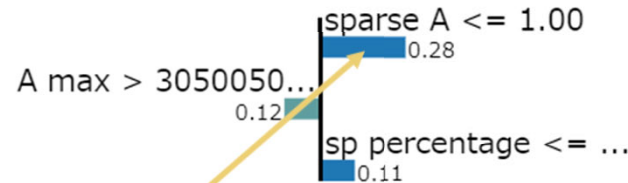
1dgemmを97%予測

Prediction probabilities



NOT 1dgemm

1dgemm



Feature

Value

Feature	Value
sparse A	0.00
A max 99999999999999991433150857216.00	0.00
sp percentage	50.45

[1.00000e+03 5.04484e+01 1.00000e+29 0.00000e+00 0.00000e+00 4.00000e+00

- 全行列が密行列の時の計算: CPUのdgemmか、GPUのdgemmのみ
- 今回の学習データ: CPUのdgemmが速い

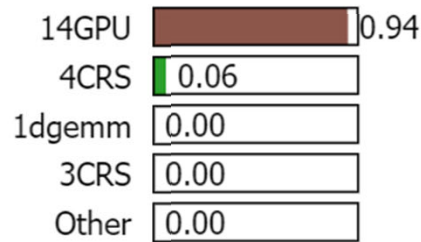
→ 分解行列の疎行列の数 (sparse A) が0の場合、
1dgemmを予測する要因として正しい



LIMEによる分析結果(3/3)

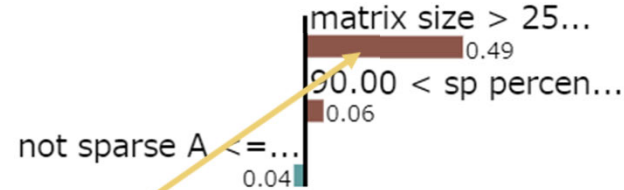
14GPU(SpMM)を94%予測

Prediction probabilities



NOT 14GPU

14GPU



Feature Value

Feature	Value
matrix size	4000.00
sp percentage	90.88
not sparse A	0.00

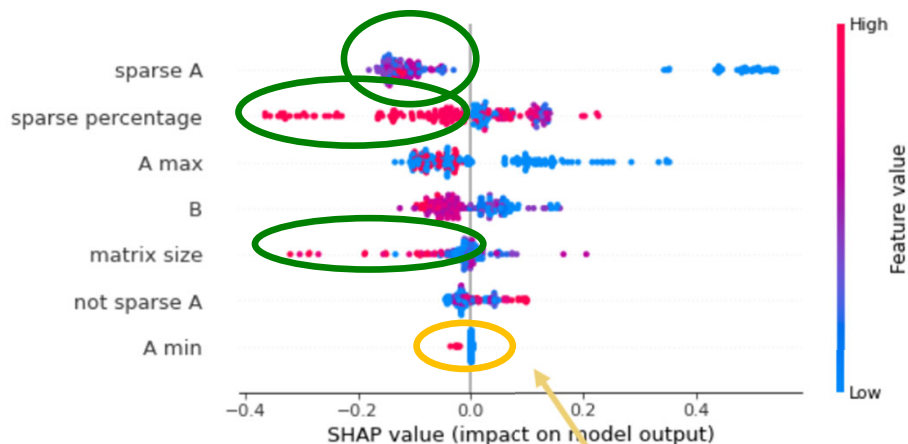
[4.0000000e+03 9.0877031e+01 1.0000000e+29 0.0000000e+00 5.0000000e+00
0.0000000e+00 6.0000000e+00]

- 学習データ: 行列サイズが大きいと14GPU (CRS形式でのSpMM(GPU))が速い傾向がある
→ (この行列: 乱数行列では)
14GPUの予測の解釈として正しい

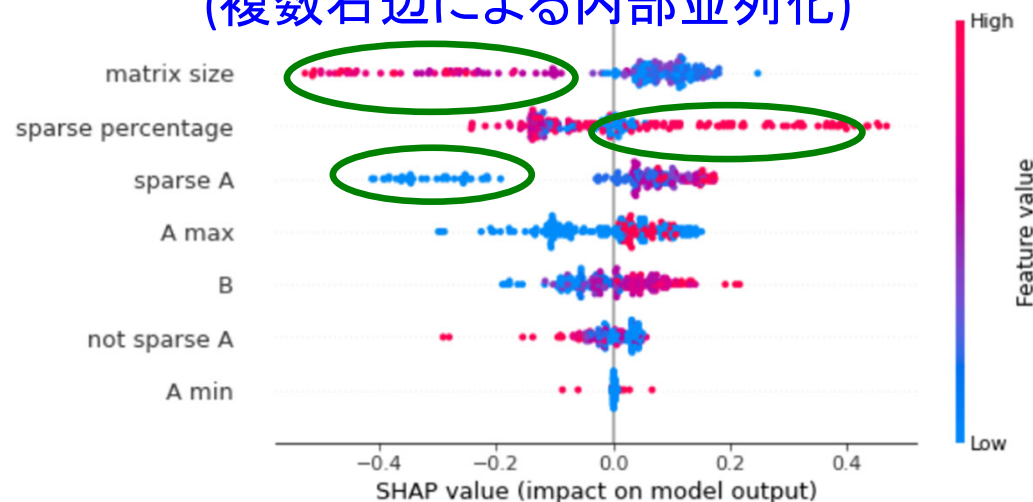
SHAPによる分析結果

- 「SHAP値の絶対値が大きく(要因が大きい)、色が近く(同じような説明変数の値)、塊が大きい、もしくは並んでいる(ケース数が多い)」箇所注目

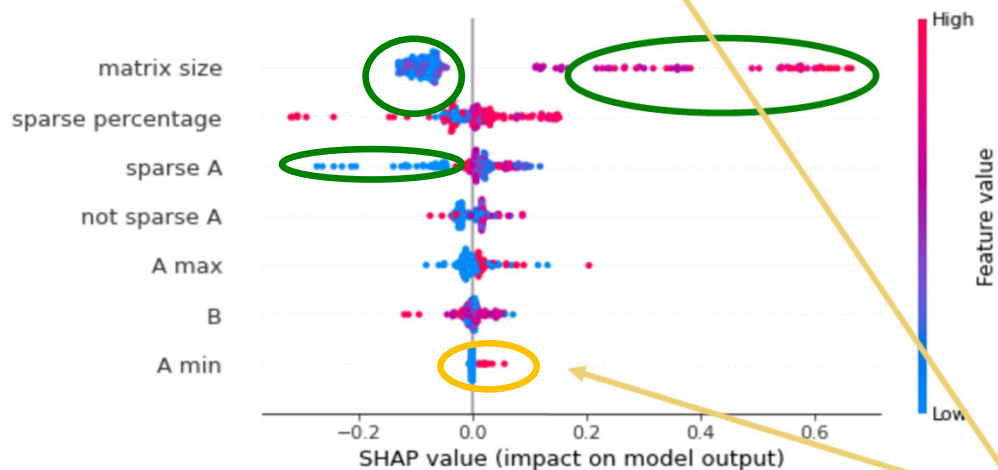
(1) 予想: dgemm



(2) 予想: CRS形式でのSpMV
(複数右辺による内部並列化)



(3) 予想: CRS形式でのSpMM(GPU)



- (1) dgemm予想: Aの疎な分解数(-)、入力行列の疎度(-:値が大)、行列サイズ(-:値が大)の要因 →妥当
- (2) CRS (SpMV) 予想: 行列サイズ(-:値が大)、入力行列の疎度(+:値が大)、Aの疎な分解数(-:値が小)の要因 →妥当
- (3) CRS (SpMM) (GPU) 予想: 行列サイズ(-:値が小、+:値が大)、Aの疎な分解数(-:値が小)の要因 →妥当
- Aの要素の最小値は、ほとんどの場合0 →予想に影響せず、良い説明変数ではない ←指摘されて分かる事項(この行列の場合には)

まとめ(XAI)

- ▶ 高精度行列積ライブラリの実装選択ATに機械学習を適用した事例で、説明可能なAI (XAI) ツールを利用しAI出力結果を分析
- ▶ 局所説明ツールLIME、大局説明ツールSHAP、ともに、対象の特性を考慮し、妥当な分析結果が得られることを確認した

今後の課題

1. XAIを活用したAT方式の開発
 - ▶ 自動で、寄与の大きな説明変数を学習を強化する
 - ▶ 寄与が少ないなら説明変数から外して、AT時間を短縮
2. XAI環境の高速化・高度化
 - ▶ 特にSHAPでフレームワーク(scikit-learn)の疎行列演算の高速化
 - ▶ 1のAT連結: 説明結果から学習を(なんども)やり直す
→機械学習のパラメタ並列実行・分散並列高速化が必須
 - ▶ →スパコンでのXAI環境の整備 →AT機構の高度化