多倍長精度基本線形計算の高速化とその応用

静岡理工科大学 幸谷智紀

kouya.tomonori@sist.ac.jp

第20回HPC-Physワークショップ 2023年12月8日(金)







- I. 多倍長精度浮動小数点計算の目的
- II. 基本線形計算の最適化・行列乗算 (xGEMM)を中心に

Ⅲ. 応用事例

- A) 混合精度反復改良法
- B) 陰的Runge-Kutta法の高速化(昔の話題)
- C) 悪条件代数方程式の求解

IV. 今後の課題

I. 多倍長精度浮動小数点計算の目的

多倍長精度数値計算の目的



なぜ多倍長精度演算 ライブラリが必要か?



2023-12-08(Fri)

ハードウェアサポートの ある標準浮動小数点数

• IEEE754 binary16(半精 度), binary32(単精度), binary64(倍精度)

ソフトウェアで実装する 多倍長精度浮動小数点数

- Double-double(DD), Triple-double(TD), Quadruple-double(QD)(,float128(binary128))
- MPFR(/GMP)







• Dekker \rightarrow double-double(DD) \rightarrow Quadruple-Double(QD)

		IEEE c	754 double——— ld_a[0]		IEEE7 d	/54 double——— d_a[1]
	±	exponent	Mantissa	±	exponent	Mantissa

Double-double(DD)

	IEEE754 double qd_a[0]			IEEE754 double qd_a[1]				IEEE754 double qd_a[2]			IEEE754 double qd_a[3]		
±	exponent	Mantissa		±	exponent	Mantissa	±	exponent	Mantissa		±	exponent	Mantissa

Quadruple-double(QD)



2023-12-08(Fri)

3. GNU MP(GMP)の高速性

GNU MP(GMP)



- 最適なアルゴリズムの選
 択
- CPUアーキテクチャごと に最適化されたアセンブ ラルーチンの採用
- T.Granlundの変質狂的な 最適化職人芸の賜物
- 複雑すぎて、mini-gmpに まとめ直す動きも

Categorization of optimization techniques for HPC

Architecture-oriented speedup technics



Optimization using computing algorithm

Long-digit real multiplication

- Karatsuba method
- Toom-Cook method
- FFT

Complex multiplication

• 3M method

Matrix multiplication

- Blocking(Tiling) algorithm
- Strassen and Winograd
- Ozaki scheme (New!)

II. 基本線形計算の最適化 ・行列乗算(xGEMM)を中心に

Optimization of multiple-precision matrix multiplication and its application

<u>xGEMM</u>		CI	ะบ		GPU			
Opt.Method	None	AVX2	OpenMP	Ozaki Scheme	CUDA	Ozaki Scheme		
DS	?	?	?	?	Mukunoki	Mukunoki		
TS	?	?	?	Utsugiri	Utsugiri	Utsugiri		
QS	?	?	?	?	?	?		
IEEE754 Binary128	MPBLAS	?	MPBLAS	Mukunoki	Mukunoki	Mukunoki		
DD	MPBLAS, BNCmatmul	Lis, MuPAT BNCmatmul	Lis,MuPAT, MPBLAS, BNCmatmul	Utsugiri	Mukunoki	Mukunoki		
TD	BNCmatmul	BNCmatmul	BNCmatmul	Utsugiri	Utsugiri	Utsugiri		
QD	MuPAT, MPBLAS, BNCmatmul	BNCmatmul	MuPAT, MPBLAS, BNCmatmul	Utsugiri	?	?		
MPFR	MPBLAS, BNCmatmul	?	MPBLAS, BNCmatmul	Utsugiri	CUMP	?		

MPLAPACK/MPBLAS https://github.com/mahonakata/mplapackMuPATYagi, H et.al(2020)Lishttps://www.ssisc.org/lis/CUMPhttps://github.com/skystar0227/CUMPBNCmatmulhttps://na-inet.jp/na/bnc/MukunokiMukunoki, D et.al.(2021)UtugiriUtsugiri(2021-2023)

 Completed implementation of optimized DD, TD, QD and MPFR matrix multiplication (xGEMM)

- Found the effectiveness of Ozaki scheme
- Confirmed the effectiveness of Ozaki scheme to LU decomposition

BNCpackからBNCmatmulへ

- BNCpack: MPFR/GMPベースの多倍長精度数値計算ライブラリ(MPIサポート)
- BNCmatmul: MPFR/GMP, DD, TD, QD精度基本線形計算の最 適化ライブラリ <u>https://github.com/tkouya/bncmatmul</u>
 - SIMD(AVX2, AVX-512)のサポート
 - OpenMPによる並列化
 - Strassen, Winograd行列乗算の実装
 - 尾崎スキームの実装
 - LU分解(直接法)
 - DLL化 + Pythonサポート

BNCmatmul's software layer



- All FP computation is supported on CPU.
- Limited functionaries but faster than MPLAPACK/MPBLAS
- Multi-component-way fxed-precision arithmetic is constructed with original ANSI C codes : Double, DD, TD, and QD
- Multi-digit-way arbitrary-precision arithmetic is based on MPFR/GMP
- Serial multiple-precision computation is runnable on Python ecosystem.

2023-12-08(Fri)

Our targeting FP precision: DD, TD, QD and MPFR



Computing environment for benchmark tests $A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times n} \implies AB \in \mathbb{R}^{n \times n}$

Each element: $(ru - 0.5) \times \exp(rn)$ *ru*: uniform rand on [0, 1], *rn*: normal dist. rand on [-1, 1]

- [Xeon] Intel Xeon W-2295 3.0GHz 18 cores, Ubuntu 20.04.3 LTS, Intel Compiler version 2021.5.0, MPLAPACK 1.0.1, MPFR 4.1.0
- [C++(icpc), C(icc)] -O3 -fp-model precise -qopenmp -axCORE-AVX2 -march=skylake -mtune=skylake -mcpu=skylake

DD, TD, QD on Xeon



2023-12-08(Fri)

MPFR 256, 512, 768bits on Xeon



第20回HPC-Physワークショップ

A quick summary of optimization techniques for multiple precision matrix multiplication



Test linear system of equations



As an application of optimized matrix multiplication, we implemented LU decomposition in the direct method for various benchmark tests, including the Top500, and measured its utility on Xeon. The corresponding n-dimensional linear system of equation is:

$$\mathbf{lx} = \mathbf{b},\tag{2}$$

where $A \in \mathbb{R}^{n \times n}$, $\mathbf{x} \in \mathbb{R}^n$, and $\mathbf{b} \in \mathbb{R}^n$ become ill-conditional as follows:

A: The diagonal matrix is $D = \text{diag}[d_1 \cdots d_n]$, where $d_i := 10^{-26(i-1)/n}$. The random matrix is $R \in \mathbb{R}^{n \times n}$; then, A is calculated as $A := RDR^{-1}$ using the mpmath Python library. Therefore, the condition number of A is $\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = 10^{26(n-1)/n}$, which requires super-DD precision arithmetic to gain accuracy.

x, **b**: After setting $\mathbf{x} = [0 \ 1 \ \cdots \ n-1]^T$, $\mathbf{b} := A\mathbf{x}$ is calculated using mpmath.

DD, TD, QD LU on Xeon



- Reduced the effectiveness due to rectangle matrix multiplication used in LU decomposition
- Trend of optimization is similar with square matrix multiplication 2023-12-08(Fri) 第20回HPC-Physワークショップ

MPFR 256, 512, 768 bits on Xeon

 LU w/ Ozaki scheme is not faster than row-wised LU decomposition over 768bitprecision MPFR arithmetic

2023-12-08(Fri)

Summary: Ozaki scheme is effective for middle-precision LU decomposition

Prec.	Method	K	Second	Rel.Err.		Prec.	Method	K	Second	Rel.Err.
DD	Rgetrf	N/A	15.8	2.4E - 01		MPFR	Rgetrf	N/A	398.6	2.1E - 50
106bits	Normal LU	1	5.1	1.9E - 01		256 bits	Normal LU	1	250.0	1.5E - 49
	Strassen+AVX2	32	7.4	1.8E + 01			Strassen	96	273.7	8.3E - 50
	OZ 6	128	4.5	9.4E + 01			OZ 13	320	167.6	3.2E - 50
TD	Normal LU	1	118.6	3.9E - 17		MPFR	Rgetrf	N/A	492.3	2.5E - 127
159bits	Strassen+AVX2	32	95.6	5.9E - 17		512 bits	Normal LU	1	341.1	2.6E - 127
	OZ 7	96	19.8	1.7E - 17			Strassen	32	390.4	3.1E - 126
QD	Rgetrf	N/A	207.4	3.8E - 34			OZ 25	576	333.4	4.1E - 126
212bits	Normal LU	1	155.3	1.4E - 33		MPFR	Rgetrf	N/A	627.8	5.1E - 205
	Strassen+AVX2	64	180.8	1.5E - 32		768 bits	Normal LU	1	481.4	1.9E - 203
	OZ 12	160	83.2	5.8E - 33			Strassen	32	491.1	3.6E - 203
L	I	1		1	ļ		OZ 38	736	536 6	3.7E - 203

Complex GEMM: 4M and 3M methods

• Multiplication of Complex matrices

 $A \in \mathbb{C}^{m \times l}, B \in \mathbb{C}^{l \times n} \implies AB \in \mathbb{C}^{m \times n}$

• 4M method : 4 real multiplications, 2 real additions

 $AB = \{\operatorname{Re}(A)\operatorname{Re}(B) - \operatorname{Im}(A)\operatorname{Im}(B)\} + \{\operatorname{Re}(A)\operatorname{Im}(B) - \operatorname{Im}(A)\operatorname{Re}(B)\} \cdot i \in \mathbb{C}^{m \times n}$

• 3M method : 3 real multiplications, 5 real additions

$$\begin{split} T_1 &= \operatorname{Re}(A)\operatorname{Re}(B) \in \mathbb{R}^{m \times n}, T_2 = \operatorname{Im}(a)\operatorname{Im}(b) \in \mathbb{R}^{m \times n} \\ AB &= (T_1 - T_2) + \{\left(\operatorname{Re}(A) + \operatorname{Im}(A)\right)\left(\operatorname{Re}(B) + \operatorname{Im}(B)\right) - T_1 - T_2\} \cdot \mathbf{i} \in \mathbb{C}^{m \times n} \end{split}$$

Implementation of 3M method

1. MPC library <u>https://www.multipleprecision.org/mpc/</u>

- 2. Bini: MPsolver GNU MP's mpf-based 3M method
- 3. BLIS: 3M and 4M CGEMM based on DGEMM $% \left(A_{1}^{2}\right) =0$

Benchmark test of CGEMM(1/3)

$$A, B \in \mathbb{C}^{n \times n} \qquad \qquad \Box \Longrightarrow \qquad C := AB$$

Real and imaginary parts of each element $(ru - 0.5) \times \exp(rn)$

[CPU and OS] Intel Xeon W-2295 3.0GHz 18 cores, Ubuntu 20.04.3 LTS [C/C++] Intel Compiler version 2021.5.0 [Compiler Option] -O3 -std=c++11 -fp-model precise [with AVX2] -axCORE-AVX2 -march=skylake -mtune=skylake mcpu=skylake [MPLAPACK] 2.0.1, GNU MP 6.2.1, MPFR 4.1.0, MPC 1.2.1

Benchmark test of CGEMM(2/3): DD and QD prec.

- Ozaki scheme(DD, QD + IMKL DGEMM)
 - DD prec. \rightarrow 6 divisions
 - QD prec. \rightarrow 12 divisions
- There is no difference of accuracy between 4M and 3M methods

Benchmark test of CGEMM(3/3): 256bits and 768bits

Complex Matrix Multiplication: MPC 256bits, Xeon

 MPBLAS
 OZ_3M 12
 OZ_3M 13

 OZ_3M 14
 OZ_4M 12
 OZ_4M 13
 OZ_4M 14

Complex Matrix Multiplication: MPC 768bits, Xeon

 MPBLAS
 ----- Strassen
 OZ_3M 36
 OZ_3M 37

 OZ_3M 38
 OZ_4M 36
 OZ_4M 37
 OZ_4M 38

• MPBLAS and Strassen are built on MPC

- 3M Ozaki scheme is made of IMKL DGEMM
- 256bits→13 divisions
 →Ozaki scheme is the fastest
- 768bits→37 divisions→Strassen is the fastest
- The same result as on real GEMM benchmark

2023-12-08(Fri)

第20回HPC-Physワークショップ

Complex LU : MPC 256, 512, 768 bits on Xeon

2023-12-08(Fri)

第20回HPC-Physワークショップ

Summary: Ozaki scheme is effective for middle-precision LU decomposition

TABLE II: Minimum computational time (s) of MPC LUTABLE III: Minimum computational time (s) of MPC LUdecomposition and maximum relative errors of x on Xeondecomposition and maximum relative errors of x on EPYC

Prec.	Method	K	Second	Rel.Err.
MPC	Cgetrf	N/A	188.1	2.9E - 71
256 bits	Normal LU	1	134.1	4.5E - 70
	Strassen	992	136.1	1.9E - 70
	OZ 13	352	116.6	6.9E - 70
MPC	Cgetrf	N/A	236.7	6.4E - 149
512 bits	Normal LU	1	181.7	3.7E - 147
	Strassen	992	187.6	8.0E - 147
	OZ 25	992	209.0	3.4E - 147
MPC	Cgetrf	N/A	313.7	1.6E - 225
768 bits	Normal LU	1	253.0	2.6E - 224
	Strassen	992	257.5	1.1E - 224
	OZ 37	992	287.6	1.5E - 224

Prec.	Method	K	Second	Rel.Err.
MPC	Cgetrf	N/A	275.9	2.9E - 71
256 bits	Normal LU	1	184.1	4.5E - 70
	Strassen	32	182.8	1.3E - 70
	OZ 13	384	171.1	3.1E - 70
MPC	Cgetrf	N/A	327.7	6.4E - 149
512 bits	Normal LU	1	237.8	3.7E - 147
	Strassen	64	240.7	3.2E - 147
	OZ 25	992	279.6	3.4E - 147
MPC	Cgetrf	N/A	428.4	1.6E - 225
768 bits	Normal LU	1	334.9	2.6E - 224
	Strassen	32	333.5	1.6E - 224
	OZ 37	992	377.5	1.5E - 224

- The computational time is slower than on the Xeon: about 20% slower for MPC alone, and up to 40¥% slower for the Ozaki scheme with DGEMM.
- There is little difference in the maximum relative error. $g_{2023-12-08(Fri)}$ $g_{2023-12-08(Fri)}$

List of categorized functions in BNCmatmul

Ⅲ. 応用事例

A) 混合精度反復改良法

Mixed precision iterative refinement method is to reduce computational cost by combining short *S* digits arithmetic and long *L* digits arithmetic (S << L). (cf.) Buttari, Alfredo, et al. International Journal of High Performance Computing A pplications 21.4 (2007): 457-466

連立一次方程式: $A \mathbf{x} = \mathbf{b}, A \in \mathbb{R}^{n \times n}, \mathbf{b}, \mathbf{x} \in \mathbb{R}^{n}$

第20回HPC-Physワークショップ

2023-12-08(Fri)

混合精度反復改良法のPythonコード

59 c	<pre>lef iterative_refinement_dd_mp(a, b, maxtimes, rtol, atol):</pre>	789	# repeat iterative refinement process
60	dim = b.rows	790	for itimes in range(maxtimes):
61		791	# Compute residual in long precision
62	# Initialize	792	res = b - a * x
63	af_ch = (ct.c_long * (dim))()	793	
64	af = init_ddmatrix(dim, dim)	794	# until r_i _2 < sqrt(n) * reps * A _F * x_i _2
65	<pre>bf = init_ddvector(dim)</pre>	795	<pre>norm_x = mpmath.norm(x); norm_res = mpmath.norm(res)</pre>
66	<pre>xf = init_ddvector(dim)</pre>	796	<pre>if norm_res < mpmath.sqrt(mpmath.mpf(dim)) * rtol * norm_a * norm_x +</pre>
67			atol:
68	<pre>x = init_mpvector(dim)</pre>	797	break
69	<pre>res = init_mpvector(dim)</pre>	798	
70		799	<pre># normalization: res := coef * res</pre>
71	<pre>resf = init_ddvector(dim)</pre>	800	<pre>normalization_coef = mpmath.mpmathify(1) / norm_res</pre>
72	z = init_mpvector(dim)	801	<pre>res = normalization_coef * res</pre>
73	zf = init_ddvector(dim)	802	
74		803	# Demote the residual from long precision to short precison
75	# norm_a := A _F	804	<pre>resf = get_ddvector_mpvec(resf, res)</pre>
76	norm_a = mpmath.mnorm(a, 'fro')	805	
77		806	<pre># Back-solve on short precision residual and short precision factors</pre>
78	# Make short precision copy of A and b	807	SolveDDLSPM(zf, af, resf, af_ch)
79	af = get_ddmatrix_mpmat(af, a)	808	
80	<pre>bf = get_ddvector_mpvec(bf, b)</pre>	809	# Promote the correction from short precision to long precision
81		810	<pre>z = set_ddvector_mpvec(z, zf)</pre>
82	# Compute LU factorization in short precision	811	
83	DDLUdecompPM(af, af_ch)	812	# reverse normalizationi
84	SolveDDLSPM(xf, af, bf, af_ch)	813	z = norm_res * z
85		814	
86	# Promote te solution from short precision to long precision	815	# Update solution in long precision
87	<pre>x = set_ddvector_mpvec(x, xf)</pre>	816	$\mathbf{x} = \mathbf{x} + \mathbf{z}$

- A We produce $A := RDR^{-1}$ with the appropriate precision mpmath, where the real diagonal matrix $D = \text{diag}[d_1 \cdots d_n]$ $(d_i := 10^{-26(i-1)/n})$ and real random matrix $R \in \mathbb{R}^{n \times n}$. Therefore, $\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = 10^{26(n-1)/n}$.
- \mathbf{x}, \mathbf{b} The true solution is $\mathbf{x} = [0 \ 1 \ \cdots \ n-1]^T$, and we calculate $\mathbf{b} := A\mathbf{x}$ with mpmath.

Core	i9: $n = 200$	0		EP'	YC: $n = 200$	00
Prec. iter. refinement	Time(s)	Max.RE	#Iter.			
DD+MP (424 bits)	190.7	2.2E-78	15	294.2	2.2E-78	16
DD+MP (AVX2)	177.0	2.7 E- 77	15	279.4	2.7 E- 77	15
TD+MP	171.0	1.0E-80	2	241.0	1.0E-80	2
TD+MP (AVX2)	131.3	8.8E-80	2	200.0	8.8E-80	2
QD+MP	409.8	7.5E-99	1	564.6	7.5E-99	1
QD+MP (AVX2)	176.5	1.8E-98	1	287.9	1.8E-98	1
MP direct (424 bits)	17940.8	8.8E-99	N/A	27700.3	8.8E-99	N/A

- AVX2による高速化の効果は DD+MPよりTD+MP, QD+MPで顕著(反復 回数に依存)
- 最小の計算時間はTD+MPとなった

B) 陰的Runge-Kutta法の高速化 (昔の話題)

IVP of n-th dimensional ODE to be solved

$$\begin{cases} \frac{d\mathbf{y}}{dt} = \mathbf{f}(t, \mathbf{y}) \in \mathbb{R}^n \\ \mathbf{y}(t_0) = \mathbf{y}_0 \end{cases}$$
(1)
Integration Interval: $[t_0, \alpha]$

We suppose that this above ODE has the unique solution, so Lipschize constant L > 0 exists to be satisfied such as

$$||\mathbf{f}(t,\mathbf{v}) - \mathbf{f}(t,\mathbf{w})|| \le L||\mathbf{v} - \mathbf{w}||$$
(2)

for $\forall \mathbf{v}, \mathbf{w} \in \mathbb{R}^n$, $\forall t \in [t_0, \alpha]$.

\downarrow

1D Brusselator problem and K-S eq. has large L >> 1, so they are called "stiff problems."

Skeleton of m stages IRK methods

Discretization :
$$t_0, t_1 := t_0 + h_0, ..., t_{k+1} := t_k + h_k...$$

When we calculate the approximation $\mathbf{y}_{k+1} \approx \mathbf{y}(t_{k+1})$ from the former $\mathbf{y}_k \approx \mathbf{y}(t_k)$, the following two steps are executed: (A) Inner iteration: Solve the nonlinear equation for unknown $\mathbf{Y} = [Y_1 \dots Y_m]^T \in \mathbb{R}^{mn}$. $\begin{cases} Y_1 = \mathbf{y}_k + h_k \sum_{j=1}^m a_{1j} \mathbf{f}(t_k + c_j h_k, Y_j) \\ \vdots \\ Y_m = \mathbf{y}_k + h_k \sum_{j=1}^m a_{mj} \mathbf{f}(t_k + c_j h_k, Y_j) \end{cases}$

(B) Calculate the next approximation \mathbf{y}_{k+1} with the above \mathbf{Y} .

$$\mathbf{y}_{k+1} := \mathbf{y}_k + h_k \sum_{j=1}^m b_j \mathbf{f}(t_k + c_j h_k, Y_j)$$

(3)

The whole algorithm of accelerated IRK method

Initial guess: $\mathbf{Y}_{-1} \in \mathbb{R}^{mn}$ For l = 0, 1, 2, ... Simplified Newton iteration (1) $\mathbf{Y}_l := [Y_1^{(l)} \ Y_2^{(l)} \ \dots \ Y_m^{(l)}]^T$ (2) $C := I_m \otimes I_n - h_k X \otimes J$, Compute $||C||_F$ (3) $\mathbf{d} := (W^T B \otimes I_n)(-\mathbf{F}(\mathbf{Y}_l))$ (4) Solve $C\mathbf{x}_0 = \mathbf{d}$ for \mathbf{x}_0 (S) For $\nu = 0, 1, 2, \dots$ Mixed precision iterative refinement (5) $\mathbf{r}_{\nu} := \mathbf{d} - C\mathbf{x}_{\nu}$ $| (6-1) \mathbf{r}'_{\nu} := \mathbf{r}_{\nu} / ||\mathbf{r}_{\nu}|| (S)$ (6-2) Solve $C\mathbf{z} = \mathbf{r}'_{\nu}$ for \mathbf{z} (S) (6-3) $\mathbf{x}_{\nu+1} := \mathbf{x}_{\nu} + ||\mathbf{r}_{\nu}||\mathbf{z}|$ (6-4) Check convergence $\Rightarrow \mathbf{x}_{\nu_{stop}}$ (7) $\overline{\mathbf{Y}_{l+1} := \mathbf{Y}_l + (W \otimes I_n) \mathbf{x}_{\nu_{ston}}}$ Check convergence \Rightarrow **Y**_{*l*_{ston}} $\mathbf{Y} := \mathbf{Y}_{l_{stop}} = [Y_1 \ Y_2 \ \dots \ Y_m]^T$ $\mathbf{y}_{k+1} := \mathbf{y}_k + h_k \sum_{j=1}^m b_j \mathbf{f}(t_k + c_j h_k, Y_j)$

1D Brusselator Problem

$$\begin{cases} \frac{\partial u}{\partial t} = 1 + u^2 v - 4 + 0.02 \cdot \frac{\partial^2 u}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \psi \end{cases}$$

$$\begin{cases} \frac{\partial u}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial^2 v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - u^2 v + 0.02 \cdot \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial x^2} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial t} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial t} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial t} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial t} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial t} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial t} \\ \frac{\partial v}{\partial t} = 3u - \frac{\partial v}{\partial t} \\ \frac{\partial v}{\partial t} \\ \frac{\partial v}{\partial t} \\ \frac{\partial$$

$$\begin{cases} \frac{du_i}{dt} = 1 + u_i^2 v_i - 4 + 0.02 \cdot \frac{u_{i+1} - 2u_i + u_{i-1}}{(\Delta x)^2} \\ \frac{dv_i}{dt} = 3u_i - u_i^2 v_i + 0.02 \cdot \frac{v_{i+1} - 2v_i + v_{i-1}}{(\Delta x)^2} \\ \Delta x = 1/(N+1) \\ u_0(t) = u_{N+1}(t) = 1, \ v_0(t) = v_{N+1}(t) = 3, \\ u_i(0) = 1 + \sin(2\pi i \Delta x), \ v_i(0) = 3 \\ (i = 1, 2, ..., N = 500) \\ \text{Integration Interval} : [0, 10] \end{cases}$$

- The discretization of x direction is fixed for benchmarking
- Double precision left preconditioned BiCGSTAB method in DP-MP mixed precision iterative refinement.

C) 悪条件代数方程式の求解

実係数代数方程式の固有値解法

代数方程式向け同時反復法

- 実装と並列化が簡単。
- ・解が近接するほど、計算精 2
 度を上げる必要がある。
- 固有値問題としては良条件 でも、代数方程式としては 悪条件になることが多い。

2 次解法

3次解法

 $z_i^{(k+1)} := z_i^{(k)} - \frac{q_n(z_i^{(k)})}{\prod_{j=1, j \neq i}^n (z_i^{(k)} - z_j^{(k)})}$

$$z_i^{(k+1)} := z_i^{(k)} - \frac{\frac{p_n(z_i^{(k)})}{p'_n(z_i^{(k)})}}{1 - \frac{p_n(z_i^{(k)})}{p'_n(z_i^{(k)})}} \sum_{\substack{j=1\\j \neq i}}^n (z_i^{(k)} - z_j^{(k)})^{-1}$$

 $\mathbf{z}_k = [z_1^{(k)} \ z_2^{(k)} \ \dots \ z_n^{(k)}]^T \in \mathbb{C}^n$

混合精度計算による高速化

- 1. DD精度のxGEEV(Rgeev)でコンパニオン行列Cの固有値を求める 2. これを同時反復法の初期値とし、係数以上の精度で反復する
- Abirthの初期値を比較対象として使用。

- 小澤の初期値も試みたが、高速化できるものとそうではないものがあり、DD精度固有値の初期値が最も良かった。
- 停止条件: $|z_i^{(k+1)} z_i^{(k)}| \le \varepsilon_{\text{rel}} |z_i^{(k)}| + \varepsilon_{\text{abs}}$

Wilkinsonの例題:全て実数解

$$\alpha_i = i$$
 となる $p_n(x) = \prod_{i=1}^n (x-i)$ を展開して与える

• Chebyshev積分公式の分点:ほぼ全て共役複素数解

$$a_n = 1 \quad \bigoplus \begin{cases} a_{n-(2k-1)} := 0\\ a_{n-2k} := -\sum_{2j+1}^k a_{n-2(k-j)}\\ \text{ここで, } k = 1, 2, ..., |n/2| である. \end{cases}$$

Wilkinsonの例題: n = 128

Wilkinson n=128, DD, QD, MPREAL(512bits) Rgeev

- $a_0 = 3.8 \cdots \times 10^{215}$
- 係数は512bits計算
- コンパニオン行列をDD, QD, MPREAL 512bits精度で与え てRgeevを適用
- DD, QD精度は全く精度が出ていない→初期誤差の影響
- MPREAL 512bits計算でよう やく実数解が現れる

Chebyshev積分公式の分点計算: n = 256

Chebyshev n=256, DD, QD, MPREAL(512bits), Rgeev

- 係数は512bits計算
- コンパニオン行列をDD, QD, MPREAL 512bits精度で与えて Rgeevを適用
- DD, QD精度は全く精度が出ていない→初期誤差の影響
- MPREAL 512bits計算でようやく 精度の良い複素数解が現れる

Wilkinsonの例題: n=128, 512bits計算

n=128	Wilkinson					$arepsilon_{\mathrm{rel}}$:=	= 7.5 imes 10	$)^{-145}, \varepsilon_{\mathrm{ab}}$	$_{ m os} := 1.0 \times$	10^{-300}	
EPYC	MPLAPA	CK Rgeev	DKA2(10)24bits)	DKA3(10)24bits)	DK2(1024	bits)+DD	DK3(1024b	its)+DD	
Unit:		MPFR									
second	DD	512bits	Second	#lter.	Second	#Iter.	Second	#Iter.	Second	#Iter.	
1 Thread			79.7	1374	71.1	691	30.8	538	51.5	512	
2 Threads			39.8	1374	36.6	691	15.4	535	26.1	509	
4 Threads	0.106	5 51	20	1374	18.7	691	8.14	562	14	512	
8 Threads	0.100	5.1	10.1	1374	8.96	691	1.43	195	1.42	99	
16 Threads			5.15	1374	4.55	691	2.07	557	3.32	511	
24 Threads			3.88	1374	3.74	691	1.52	541	2.69	511	
n=128	Wilkinson										
Xeon	MPLAPA	CK Rgeev	DKA2(1024bits)		DKA3(1024bits)		DK2(1024bits)+DD		DK3(1024bits)+DD		
Unit:		MPFR									
second	DD	512bits	Second	#lter.	Second	#Iter.	Second	#Iter.	Second	#Iter.	
1 Thread			58.6	1374	50.8	691	22.3	538	37.1	512	
2 Threads			29.2	1374	26.2	691	11.3	535	18.8	509	
4 Threads	0.00	2 5	15.3	1374	13.7	691	6.18	562	9.93	512	
8 Threads	0.09	3.5	7.78	1374	6.86	691	1.1	195	0.996	99	
16 Threads				4.73	1374	4.16	691	1.86	557	2.97	511
18 Threads			4.88	1374	4.37	691	1.9	557	3.05	511	

赤字部分はMPFR Rgeev計算時間より高速であることを示す。

000

Wilkinsonの例題: EPYC vs. Xeon

- MPFR Rgeevと同時反復法は同程度の相対誤差の近似値を得ている。
- 8 Threadsで効率が急上昇しているのは, DD精度の近似値を用いた場合の 反復回数が激減したことによる。

Chebyshev積分公式の分点: n=256, 256 bits係数

 $\varepsilon_{\rm rel} := 8.6 \times 10^{-68}, \ \varepsilon_{\rm abs} := 1.0 \times 10^{-300}$

11 230	encoyoner			,						
EPYC	MPLAPA	CK Rgeev	DKA2(5	12bits)	DKA3(5	12bits)	DK2(512	bits)+DD	DK3(512	bits)+DD
Unit:		MPFR								
second	DD	256bits	Second	#lter.	Second	#Iter.	Second	#Iter.	Second	#Iter.
1 Thread			203.0	1344	183	671	47	316	70.3	264
2 Threads			100.0	1344	93.8	671	7.11	96	10	73
4 Threads	1 20	FTC	51.1	1344	46.3	671	3.62	96	5.08	71
8 Threads	1.29	57.0	25.3	1344	23.1	671	5.77	308	8.92	264
16 Threads			13	1344	11.7	671	0.918	95	1.22	71
24 Threads			8.96	1344	8.1	671	2.02	263	3.1	263

Chebyshev n=256, 256bits-Rgeev, 512bits-

DK

Chehyshev

n=256

- Rgeevと同時反復法の近似解はほぼ同じ相対誤差の近似解を 得ている。
- DK2-DD, DK3-DD(DD Rgeevを初期値とする)は反復回数の 変化が大きい。
- 赤字部分はMPFR Rgeev計算時間より高速であることを示す。

Chebyshev積分公式の分点: n=512, 512bits係数

n=512	Chebyshe	ev					$\varepsilon_{\mathrm{rel}} := i$	$.5 \times 10^{-110}$	$, \varepsilon_{\rm abs} := 1.$	0×10^{-000}
EPYC	MPLAP	ACK Rgeev	DKA2(1024bits)		DKA3(1024bits)		DK2(1024bits)+DD		DK3(1024bits)+DD	
		MPFR								
Unit: second	DD	512bits	Second	#Iter.	Second	#Iter.	Second	#Iter.	Second	#Iter.
1 Thread			2850.0	3032	2460	1510	500	552	577	357
2 Threads	DD 9.84		1370.0	3032	1210	1510	90.9	200	67.8	82
4 Threads	0.04		689	3032	614	1510	46.1	202	34.5	83
8 Threads	9.84	0.080	344	3032	307	1510	62.8	550	77.3	363
16 Threads			173	3032	174	1510	11.8	205	8.59	84
24 Threads			120	3032	106	1510	8.47	213	6.57	83

Chebyshev n=512, 512bits-Rgeev,

- Rgeevと同時反復法の近似解はほぼ同じ相対誤差の近似解を 得ている。
- DK2-DD, DK3-DD(DD Rgeevを初期値とする)は反復回数の 変化が大きい。
- 赤字部分はMPFR Rgeev計算時間より高速であることを示す。

10 - 300

Chebyshev積分公式の分点:Xeon

n=256

Xeon	MPLAPA	CK Rgeev	DKA2	512bits)	DKA3(5	12bits)	DK2(512	bits)+DD	DK3(512	bits)+DD
		MPFR								
Unit: second	DD	256bits	Second	#Iter.	Second	#Iter.	Second	#Iter.	Second	#Iter.
1 Thread			146.0	1344	127	671	32.1	316	49.2	264
2 Threads			69.0	1344	64.8	671	4.95	96	6.91	73
4 Threads	1 1 2	20 /	36.3	1344	34.2	671	2.61	96	3.59	71
8 Threads		38.4	18.6	1344	17.7	671	4.2	308	6.63	264
16 Threads			10.9	1344	10	671	0.779	95	1.07	71
18 Threads			10.8	1344	9.79	671	2.32	293	3.62	250

n=512

Xeon	MPLAPACK Rgeev		DKA2(1024bits)		DKA3(1024bits)		DK2(1024bits)+DD		DK3(1024bits)+DD	
		MPFR								
Unit: second	DD	512bits	Second	#lter.	Second	#Iter.	Second	#Iter.	Second	#Iter.
1 Thread	9.25	441.0	2030.0	3032	1740	1510	362	552	407	357
2 Threads			1020.0	3032	937	1510	67.3	200	48.2	82
4 Threads			528	3032	471	1510	35.3	202	25.5	83
8 Threads			267	3032	238	1510	48.5	550	56.5	363
16 Threads			159	3032	142	1510	10.8	205	7.85	84
18 Threads			151	3032	132	1510	26.9	539	31.7	365

IV. 今後の課題

やり残しを片付ける&ライブラリ整備

• やり残しを片付ける

- ・疎行列(SpMV)実装→Lanczos法の高速化
- BNCmatmulベースの陰的Runge-Kutta法の再実装
- 混合精度代数方程式の求解→行列の固有値問題の高速化
- ライブラリ整備
 - WASMを用いたWeb上の多倍長精度計算の高速化
 - Pythonモジュールの整備
 - ドキュメント整理