

生成 AI の活用による HPC コード GPU 化の展望

棕木 大地（名古屋大学 情報基盤センター）

2026 年 1 月 21 日

「次世代計算基盤を見据えたソフトウェア環境整備とそれを担う人材の育成に関する提言」
についての意見交換会

生成 AI 技術の現状

- ChatGPT に代表される「博士号レベルの知識」を持った AI が普及
- 長時間自律動作しツールを使いこなす **コーディング AI エージェント**
 - ▶ Claude Code (Anthropic), Codex (OpenAI), Gemini CLI (Google)
- **Vibe Coding** – AI とのカジュアルな対話を通してコーディングする開発スタイル

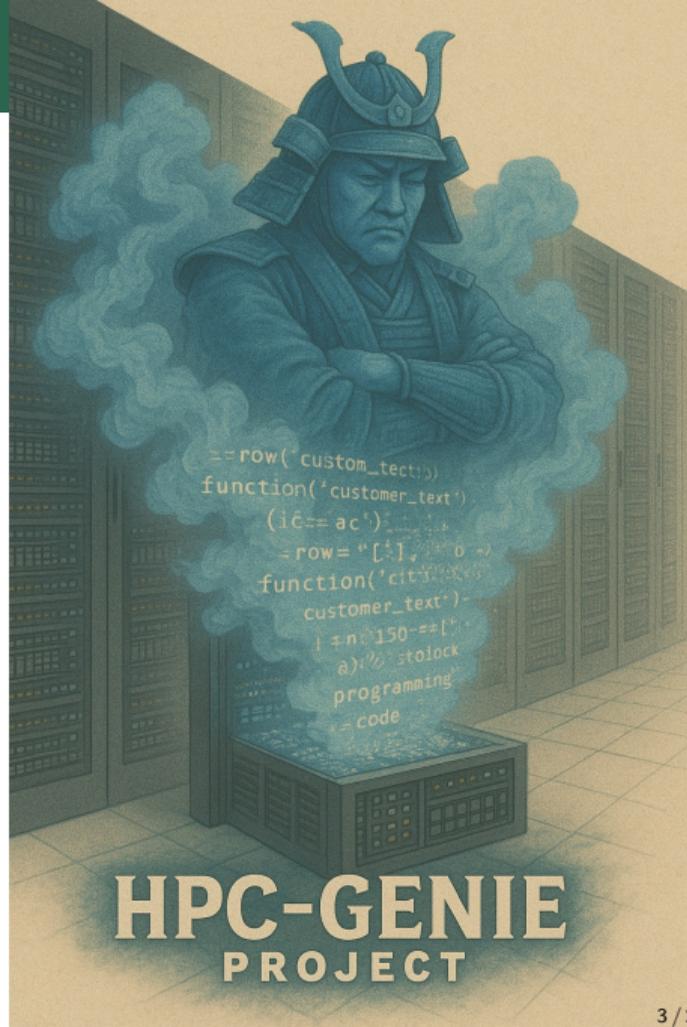
HPC 業界におけるコード生成 AI への期待

- 研究開発の効率化と人材不足の解消
- HPC の研究・開発に求められる高度な知識・技術の補完
 - ▶ 多様かつ特殊なアーキテクチャ (GPU, 分散並列, FPGA etc.), 言語 (C/C++/Fortran etc.)
 - ▶ 数値計算特有の問題 (最適なアルゴリズム選択, 計算誤差への対処 etc.)
- 世界中で研究が進んでいるがまだ黎明期か
 - ▶ ChatHPC (Oak Ridge National Laboratory) ¹⁾ – HPC 開発・研究に特化した AI アシスタント生成フレームワーク

¹⁾P. V. Lara, A. Young, J. S. Vetter, Z. Jin, S. Pophale, M. A. H. Monil, K. Teranishi, and W. F. Godoy. 2025. ChatHPC: Building the Foundations for a Productive and Trustworthy AI-Assisted HPC Ecosystem. SC '25.

HPC-GENIE

- High-Performance Computing with Generative Neural Intelligence for Execution
- 名古屋大学情報基盤センターで発足した AI による HPC コード生成研究プロジェクト (2025/4-)
- https://www.hpc.itc.nagoya-u.ac.jp/menu/hpc_genie.html
- 周回遅れ感は否めないが今から追いつき，追い越す
- 学生が活躍 – 学生の AI への関心は非常に高い



主要ミッション

- **エージェント開発**
 - ▶ マルチエージェント等の工夫により既存モデルを有効に活用
 - ▶ 資金と資源が必要な大規模モデル開発とは別のアプローチを模索
- 小規模モデル開発
 - ▶ 加えて HPC のための学習データセット構築，ベンチマークの構築
- **ローカル LLM** – 商用サービス依存を減らす
 - ▶ 課金爆発: API 利用料・サブスク料金
 - ▶ ブラックボックス: クローズドソース，研究開発のために手を入れられない
 - ▶ セキュリティ懸念: 外部サーバでコードや研究データが処理される
 - ▶ 国産技術開発
- **富岳 NEXT のための Fortran to GPU**

GPT-4.1 & o4-mini を用いた BLAS コード生成²⁾

正しく動作するコード数 (10 回生成中)

Level	Routine	非最適化コード		最適化コード	
		GPT-4.1	o4-mini	GPT-4.1	o4-mini
1	dasum	10	10	9	7
	daxpy	10	10	10	10
	ddot	10	10	8	9
	idamax	3	10	2	9
	dnrm2	10	10	7	5
	drot	10	10	8	9
2	dgemv	8	9	3	6
	dger	10	10	7	7
	dsymv	8	8	0	2
	dsyr2	8	9	2	7
	dtrmv	3	4	0	5
	dtrsv	8	9	0	4
3	dgemm	10	10	3	0
	dsymm	4	8	3	3
	dsyrk	3	10	0	1
	dsyr2k	3	10	0	0
	dtrmm	0	1	0	0
	dtrsm	0	0	0	0

- OpenAI GPT-4.1 & o4-mini – エントリーレベルの汎用 LLM
- 単純なプロンプトによる一発生成：
「BLAS の #ROUTINE#ルーチンを C 言語で実装せよ」（実際は英語）
- ルーチン名だけから動作するコードが生成される
- 生成コードは Fortran リファレンスコードと異なるコード構造が多い（なるべく短いコードを生成）
- 性能最適化コードの一発生成は不十分

²⁾D. Mukunoki, S. Hayashi, T. Hoshino, T. Katagiri, “Performance Evaluation of General Purpose Large Language Models for Basic Linear Algebra Subprograms Code Generation”, arXiv preprint arXiv:2507.04697, 2025.

椋木大地, 林俊一郎, 星野哲也, 片桐孝洋, BLAS コードを題材とした GPT モデルによる数値計算コード実装支援に関する考察, 第 200 回ハイパフォーマンスコンピューティング研究発表会 (SWoPP2025), 2025 年 8 月 4 日.

Intel GPU 向け SYCL コード生成能力の調査⁴⁾

Performance (Level 2: GB/s, Level 3: GFlops/s, Intel Arc B580, oneAPI 2025.2.1)

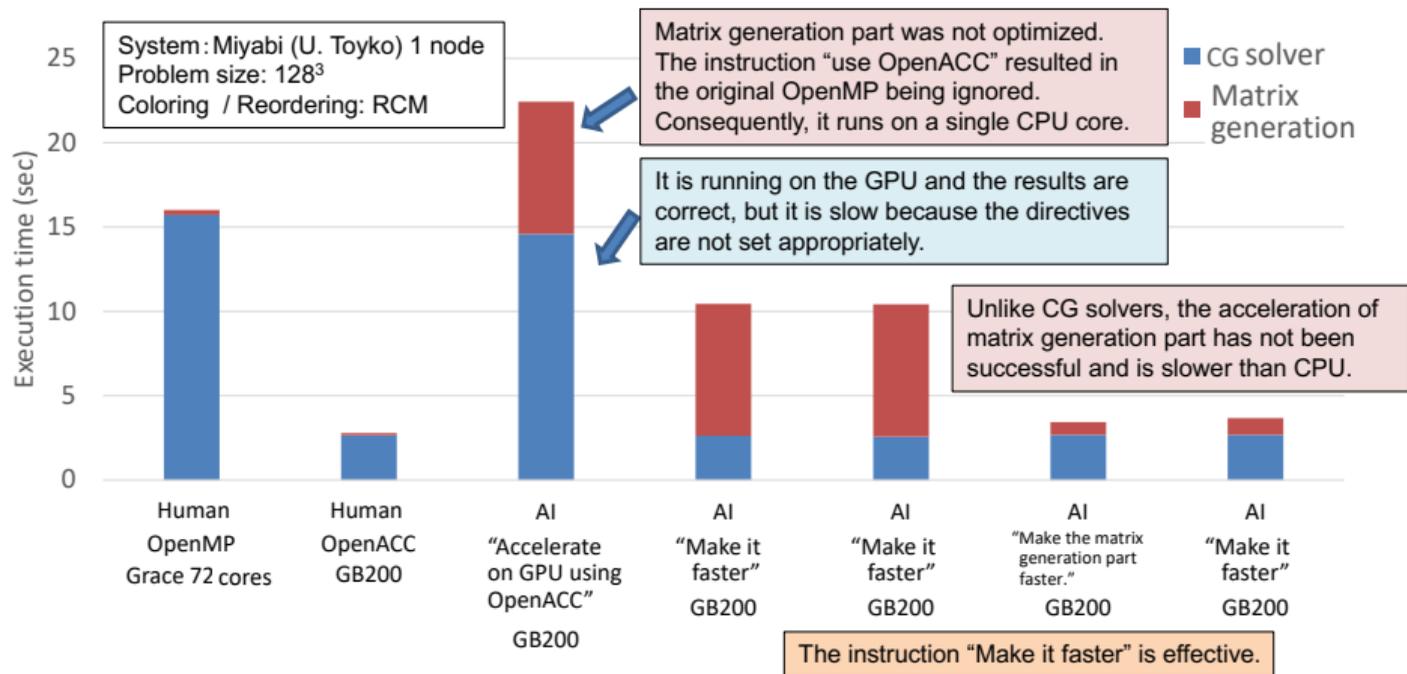
Level	Routine	MKL	run 1	run 2	run 3
2	SGEMV	413.96	408.49	328.53	316.64
	SSYMV	123.01	158.93	151.03	87.72
	STRSV	86.07	11.83	59.46	11.75
3	SGEMM	14030.35	3340.28	4380.32	3643.29
	SSYMM	8863.38	3247.80	4207.13	2359.42
	SSYRK	7613.51	3287.25	5343.23	3680.24

- Claude Code³⁾ による Intel GPU 向け SYCL コード生成 (対象: BLAS)
- 入力: CPU 向けの最適化されていないコード, 最適化手順書 (ターゲットハードウェア, コンパイル方法, テスト方法)
- プロンプト: 「手順書に従って最適化せよ」 → 「もっと最適化しろ」 → 「もっと最適化しろ」
- NVIDIA 以外の GPU に対応することで**ベンダーロックインの回避に期待**

³⁾モデル: claude-sonnet-4-5-20250929

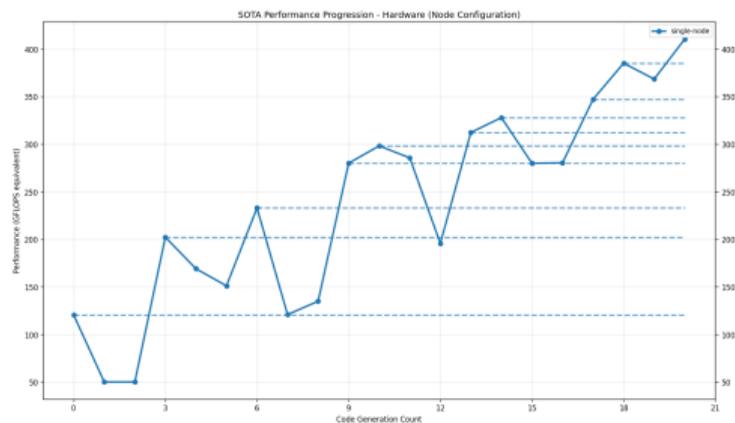
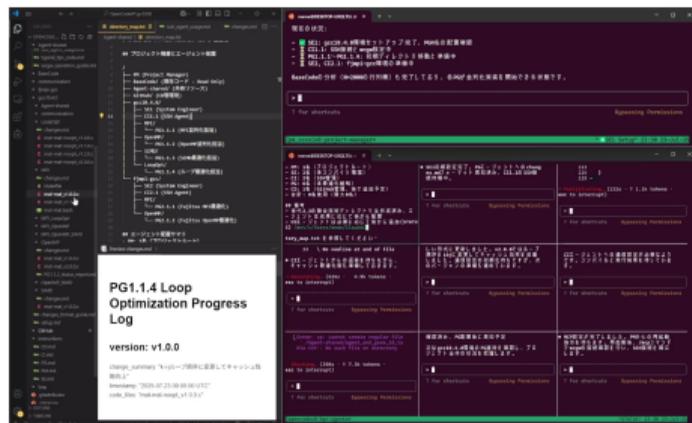
⁴⁾K. Morita, D. Mukunoki, T. Hoshino, T. Katagiri, Evaluation of the Capability of Coding AI in Generating SYCL-Based Numerical Computation Codes for Intel GPUs, SCA/HPCAsia2026, poster session, Jan. 2026 (accepted).

Claude Code による Fortran コード (GeoFEM) の GPU 化⁵⁾



GPU acceleration of GeoFEM (matrix generation + CG solver) on NVIDIA GB200

⁵⁾T. Hoshino, S. Hayashi, D. Mukunoki, T. Katagiri and T. Hanawa, Evaluating Claude Code's Coding and Test Automation for GPU Acceleration of a Legacy Fortran Application: A GeoFEM Case Study, LLM4HPCAsia 2026, 2026 (accepted).



VibeCodeHPC (主要開発者：林俊一郎 (名古屋大学 修士1年))

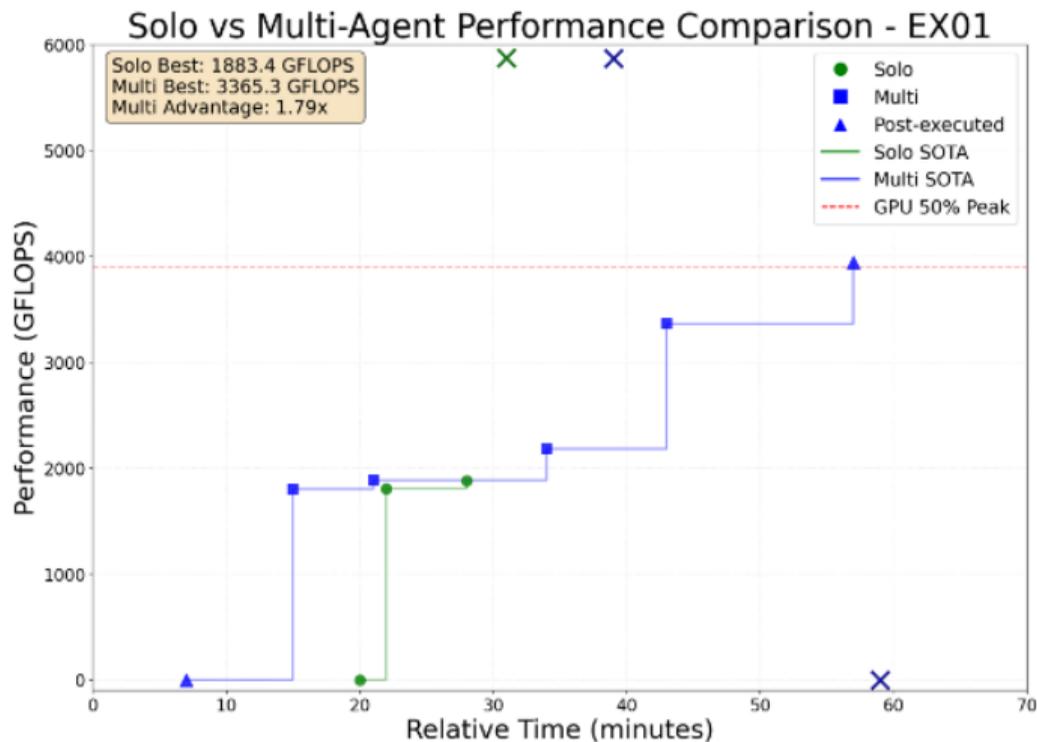
- CLI コーディングエージェント (Claude Code) を利用した**スパコン向けコード開発エージェントシステム**
- 役割の異なる複数のエージェントが連携して作業する**マルチエージェントシステム**
- Vibe Coding も可能であるが**長時間完全自律動作も可能**

⁶⁾<https://github.com/Katagiri-Hoshino-Lab/VibeCodeHPC-jp>

⁷⁾Shun-ichiro Hayashi, Koki Morita, Daichi Mukunoki, Tetsuya Hoshino, Takahiro Katagiri, VibeCodeHPC: An Agent-Based Iterative Prompting Auto-Tuner for HPC Code Generation Using LLMs, arXiv preprint arXiv:2510.00031,2025.

林俊一郎, 森田光貴, 椋木大地, 星野哲也, 片桐孝洋. VibeCodeHPC: 自動コード最適化 CLI 型マルチエージェント, 第 255 回システム・アーキテクチャ・第 202 回ハイパフォーマンスコンピューティング合同研究発表会, 2025 年 12 月 15 日.

VibeCodeHPC – ソロエージェント vs マルチエージェント



マルチエージェントにより開発失敗や要件定義違反の発生を防ぐ

DGEMM (行列積)

問題サイズ: $m=n=k=2048$, in GFlops/s

PG role = dir. config.	Agent config.	run 1	run 2	run 3
OpenMP	Multi	197.7	213.0	189.6
	Solo	–	–	–
MPI	Multi	6.7	30.7	–
	Solo	–	–	–
CUDA	Multi	1384.6	3677.2	916.2
	Solo	804.5	606.9	1097.7
OpenACC	Multi	249.8	–	–
	Solo	–	–	–
MPI+CUDA	Multi	–	–	332.1
	Solo	–	–	–
SIMD	Multi	–	183.5	–
	Solo	–	–	–

姫野ベンチマーク (ポアソンソルバー)

問題サイズ: **SMALL**, in GFlops/s

PG role = dir. config.	Agent config.	run 1	run 2	run 3
OpenMP	Multi	47.2	77.8	67.1
	Solo	21.5	23.9	39.9
MPI	Multi	40.1	57.6	61.1
	Solo	–	–	–
CUDA	Multi	21.1	24.5	23.3
	Solo	14.1	–	19.5
OpenACC	Multi	24.4	20.6	20.3
	Solo	20.5	–	–
OpenMP+MPI	Multi	68.8	54.2	–
	Solo	–	–	–

CPU: Xeon Gold 6230 (20 cores) × 2 (2.688 TFlops/s per CPU, 281.5 GB/s)

GPU: Tesla V100 × 4, (7.8 TFlops/s per GPU, 900 GB/s)

PG エージェントの作業分担は専門 (利用する技術) ごとに自動で設定される。
マルチエージェントが優位 (ソロエージェントは全部を 1 人でこなしている)。

ローカル PC だけで動作するコード自動最適化システムの試作⁸⁾

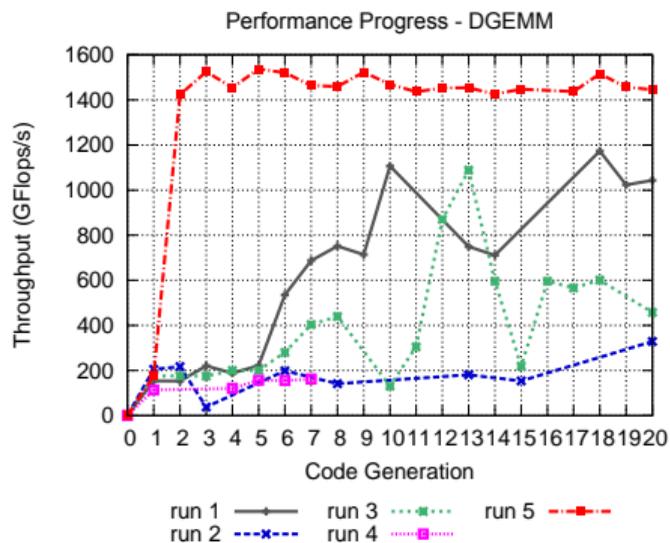
概要

- 研究開発用にコード性能最適化システムをいちから実装したい
 - ▶ 現時点ではトイ的な非常にシンプルなもの
 - ▶ 単一コードの最適化だけを反復的に繰り返す
- LLM も含めてローカル PC で動作させたい
 - ▶ **コンシューマ PC で動くレベルの LLM は 7 ヶ月前の商用最高モデル相当の性能を達成する** (Epoch AI)
 - ▶ **gpt-oss-120b** (OpenAI) を利用 – o4-mini 相当 (一世代前の ChatGPT ぐらい), Ryzen AI Max+ や MacBook Pro の 128GB メモリで動作
- 疑似マルチエージェント構成で反復的に最適化
 1. Project Manager (PM) LLM が 5 つ Programmer (PG) に対する最適化戦略を生成
 2. PG が 5 つのコードを生成
 3. 一番良いコードを採用し PM は有効・無効な最適化方法を分析
 - ▶ 以上を 20 回反復

⁸⁾ 椋木大地, 森田光貴, 林俊一郎, 三笠諒, 星野哲也, 片桐孝洋, gpt-oss-120b を用いたコード自動最適化マルチエージェントシステムの試作, 第 255 回システム・アーキテクチャ・第 202 回ハイパフォーマンスコンピューティング合同研究発表会, 2025 年 12 月 15 日.

Input Code

```
1 void exp1_opt(int N,  
2             double *A,  
3             double *B,  
4             double *C) {  
5     for (int i = 0; i < N; i++) {  
6         for (int j = 0; j < N; j++) {  
7             double sum = 0.0;  
8             for (int k = 0; k < N; k++) {  
9                 sum += A[i*N+k] * B[k*N+j];  
10            }  
11            C[i*N+j] = sum;  
12        }  
13    }  
14 }
```



- Intel Xeon Gold 6230 (20 コア) × 2, $n = 1024$ の正方行列
- 最速コードは **MKL を呼んでいた** → LLM がコードが行列積であることを認識
- 最速のフルスクラッチコード：1173.42 GFlops/s (MKL の 76%)
- 実験毎に結果が大きく異なるが、試行回数を増やせば良い結果が得られる

DGEMM (CPU) – 生成された最速コード (1/3)

```
1  static inline void pack_B_panel(const double *B_src,
2  double *B_pack, int N, int kp, int jp, int k_max,
3  int blk_width) {
4  for (int k = kp; k < k_max; ++k) {
5  const double *src = &B_src[k * N + jp];
6  double *dst = &B_pack[(k - kp) * blk_width];
7  int j = 0;
8  __mm_prefetch((const char *)src, _MM_HINT_T0);
9  for (; j + 8 <= blk_width; j += 8) {
10 __m512d v = __mm512_loadu_pd(&src[j]);
11 __mm512_storeu_pd(&dst[j], v);
12 }
13 for (; j < blk_width; ++j) { dst[j] = src[j]; }
14 }
15 void expl_opt(int N, double *restrict A, double *
16 restrict B, double *restrict C) {
17 const int BLOCK = 64;
18 omp_set_num_threads(omp_get_max_threads());
19 #pragma omp parallel for collapse(2) schedule(static)
20 for (int ii = 0; ii < N; ii += BLOCK) {
21 for (int jj = 0; jj < N; jj += BLOCK) {
22 int i_max = (ii + BLOCK > N) ? N : ii + BLOCK;
23 int j_max = (jj + BLOCK > N) ? N : jj + BLOCK;
24 int blk_width = j_max - jj;
25 int vec_width = blk_width & ~31;
26 double C_local[BLOCK * BLOCK] __attribute__((
27 aligned(64)));
28 double *C_local_aligned = (double *)__
29 builtin_assume_aligned(C_local, 64);
30 for (int i = 0; i < BLOCK * BLOCK; ++i)
31 C_local_aligned[i] = 0.0;
32 for (int kk = 0; kk < N; kk += BLOCK) {
33 int k_max = (kk + BLOCK > N) ? N : kk +
34 BLOCK;
35 double B_pack_local[BLOCK * BLOCK] __
36 attribute__((aligned(64)));
37 pack_B_panel(B, B_pack_local, N, kk, jj,
38 k_max, blk_width);
39 for (int i = ii; i + 3 < i_max; i += 4) {
40 const double *a_row0 = &A[(i + 0) * N +
41 kk];
```

```
42 const double *a_row1 = &A[(i + 1) * N +
43 kk];
44 const double *a_row2 = &A[(i + 2) * N +
45 kk];
46 const double *a_row3 = &A[(i + 3) * N +
47 kk];
48 double *c_loc0 = &C_local_aligned[(i - ii
49 ) * BLOCK];
50 double *c_loc1 = &C_local_aligned[(i - ii
51 + 1) * BLOCK];
52 double *c_loc2 = &C_local_aligned[(i - ii
53 + 2) * BLOCK];
54 double *c_loc3 = &C_local_aligned[(i - ii
55 + 3) * BLOCK];
56 for (int j = 0; j < vec_width; j += 16) {
57 __m512d c00 = __mm512_loadu_pd(&c_loc0[
58 j + 0]);
59 __m512d c01 = __mm512_loadu_pd(&c_loc0[
60 j + 8]);
61 __m512d c10 = __mm512_loadu_pd(&c_loc1[
62 j + 0]);
63 __m512d c11 = __mm512_loadu_pd(&c_loc1[
64 j + 8]);
65 __m512d c20 = __mm512_loadu_pd(&c_loc2[
66 j + 0]);
67 __m512d c21 = __mm512_loadu_pd(&c_loc2[
68 j + 8]);
69 __m512d c30 = __mm512_loadu_pd(&c_loc3[
70 j + 0]);
71 __m512d c31 = __mm512_loadu_pd(&c_loc3[
72 j + 8]);
73 for (int k = kk; k < k_max; ++k) {
74 const int k_off = k - kk;
75 __m512d b0 = __mm512_load_pd(&
76 B_pack_local[k_off * blk_width
77 + j + 0]);
78 __m512d b1 = __mm512_load_pd(&
79 B_pack_local[k_off * blk_width
80 + j + 8]);
81 __m512d a0 = __mm512_set1_pd(a_row0[
82 k_off]);
83 __m512d a1 = __mm512_set1_pd(a_row1[
84 k_off]);
```

DGEMM (CPU) – 生成された最速コード (2/3)

```
56     __m512d a2 = __mm512_set1_pd(a_row2[
57         k_off]);
58     __m512d a3 = __mm512_set1_pd(a_row3[
59         k_off]);
60     c00 = __mm512_fmadd_pd(a0, b0, c00);
61     c01 = __mm512_fmadd_pd(a0, b1,
62         c01);
63     c10 = __mm512_fmadd_pd(a1, b0, c10);
64     c11 = __mm512_fmadd_pd(a1, b1,
65         c11);
66     c20 = __mm512_fmadd_pd(a2, b0, c20);
67     c21 = __mm512_fmadd_pd(a2, b1,
68         c21);
69     c30 = __mm512_fmadd_pd(a3, b0, c30);
70     c31 = __mm512_fmadd_pd(a3, b1,
71         c31);
72 }
73     __mm512_storeu_pd(&c_loc0[j + 0], c00);
74     __mm512_storeu_pd(&c_loc0[j + 8],
75         c01);
76     __mm512_storeu_pd(&c_loc1[j + 0], c10);
77     __mm512_storeu_pd(&c_loc1[j + 8],
78         c11);
79     __mm512_storeu_pd(&c_loc2[j + 0], c20);
80     __mm512_storeu_pd(&c_loc2[j + 8],
81         c21);
82     __mm512_storeu_pd(&c_loc3[j + 0], c30);
83     __mm512_storeu_pd(&c_loc3[j + 8],
84         c31);
85 }
86     for (int j = vec_width; j < blk_width; ++
87         j) {
88         double s0 = c_loc0[j]; double s1 =
89             c_loc1[j];
90         double s2 = c_loc2[j]; double s3 =
91             c_loc3[j];
92         for (int k = kk; k < k_max; ++k) {
93             const int k_off = k - kk;
94             double b = B_pack_local[k_off *
95                 blk_width + j];
96             s0 += a_row0[k_off] * b; s1 +=
97                 a_row1[k_off] * b;
```

```
75         s2 += a_row2[k_off] * b; s3 +=
76             a_row3[k_off] * b;
77     }
78     c_loc0[j] = s0; c_loc1[j] = s1; c_loc2
79     [j] = s2; c_loc3[j] = s3;
80 }
81     for (int i = i_max - ((i_max - ii) % 4); i <
82         i_max; ++i) {
83         const double *a_ptr = &A[i * N + kk];
84         double *c_loc = &C_local_aligned[(i - ii)
85             * BLOCK];
86         for (int j = 0; j < vec_width; j += 8) {
87             __m512d c_vec = __mm512_loadu_pd(&c_loc
88                 [j]);
89             for (int k = kk; k < k_max; ++k) {
90                 const int k_off = k - kk;
91                 __m512d b_vec = __mm512_load_pd(&
92                     B_pack_local[k_off * blk_width
93                         + j]);
94                 __m512d a_vec = __mm512_set1_pd(
95                     a_ptr[k_off]);
96                 c_vec = __mm512_fmadd_pd(a_vec,
97                     b_vec, c_vec);
98             }
99             __mm512_storeu_pd(&c_loc[j], c_vec);
100         }
101         for (int j = vec_width; j < blk_width; ++
102             j) {
103             double sum = c_loc[j];
104             for (int k = kk; k < k_max; ++k) {
105                 const int k_off = k - kk;
106                 sum += a_ptr[k_off] * B_pack_local[
107                     k_off * blk_width + j];
108             }
109             c_loc[j] = sum;
110         }
111     }
112     for (int i = ii; i < i_max; ++i) {
113         double *c_ptr = &C[i * N + jj];
114         double *c_loc = &C_local_aligned[(i - ii) *
115             BLOCK];
```

DGEMM (CPU) – 生成された最速コード (3/3)

```
106     c_ptr = (double *) __builtin_assume_aligned(  
107         c_ptr, 64);  
108     c_loc = (double *) __builtin_assume_aligned(  
109         c_loc, 64);  
110     int j = 0;  
111     for (; j < vec_width; j += 8) {  
112         __m512d tmp = _mm512_loadu_pd(&c_loc[j]);  
113         _mm512_storeu_pd(&c_ptr[j], tmp);  
114     }  
115     for (; j < blk_width; ++j) {  
116         c_ptr[j] = c_loc[j];  
117     }  
118 }  
119 }
```

※掲載スペースの都合により実際に生成されたコードから自動的に挿入されたコメント、空行、ヘッダファイルの include を削除し、インデントや改行を一部修正した。

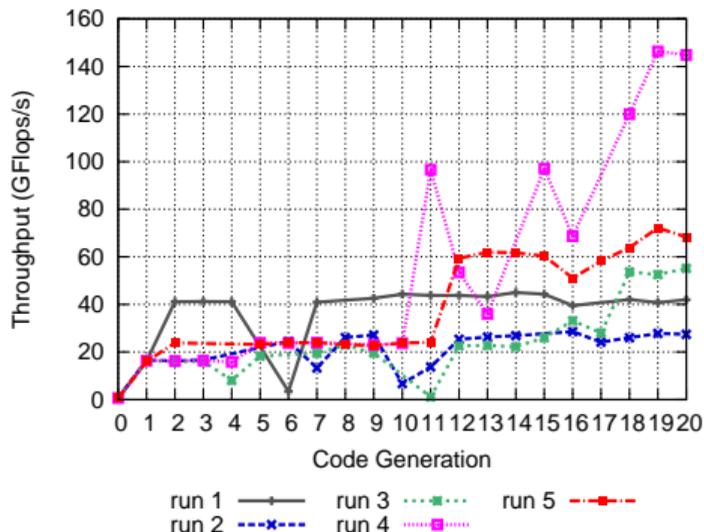
Input Code

```

1 void exp2_opt(int N,
2   double *A, double *B,
3   double *X, double alpha) {
4   for (int i = 0; i < N; i++) {
5     for (int j = 0; j < N; j++) {
6       X[i*N+j] = alpha * B[i*N+j];
7     }
8   }
9   for (int j = 0; j < N; j++) {
10    for (int i = 0; i < N; i++) {
11      for (int k = 0; k < i; k++) {
12        X[i*N+j] -= A[i*N+k]*X[k*N+j];
13      }
14      X[i*N+j] /= A[i*N+i];
15    }
16  }
17 }

```

Performance Progress - DTRSM



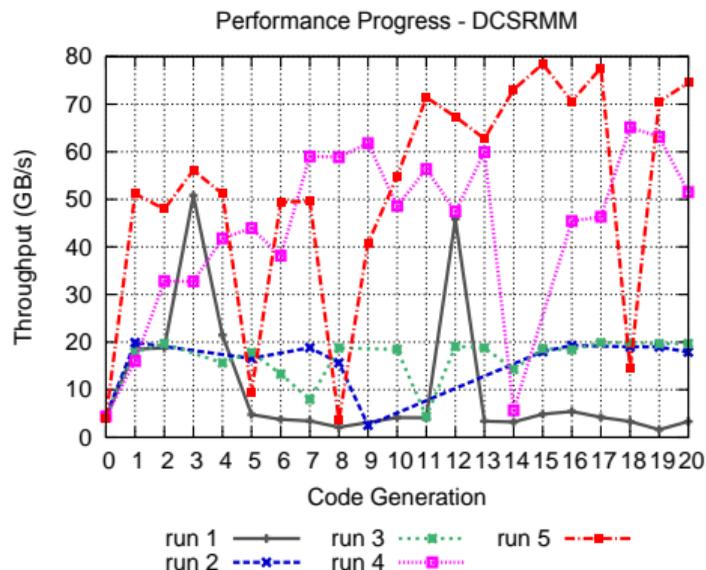
- 三角行列を係数とする方程式の計算（下三角, $n = 1024$ ） – データ依存性があり最適化の難易度が高い
- 最速：146.33 GFlops/s (run 4) – MKL (305.74 GFlops/s) の 48%

Input Code

```

1 void exp3_opt(int m, int k, int n,
2   const int *row_ptr,
3   const int *col_indices,
4   const double *values,
5   const double *B, double *C) {
6   memset(C, 0, m*n*sizeof(double));
7   for (int i = 0; i < m; i++) {
8     int row_start = row_ptr[i];
9     int row_end = row_ptr[i + 1];
10    for (int idx=row_start;
11         idx<row_end; idx++) {
12      int col = col_indices[idx];
13      double val = values[idx];
14      for (int j = 0; j < n; j++) {
15        C[i*n+j] += val*B[col*n+j];
16      }
17    }
18  }
19 }

```

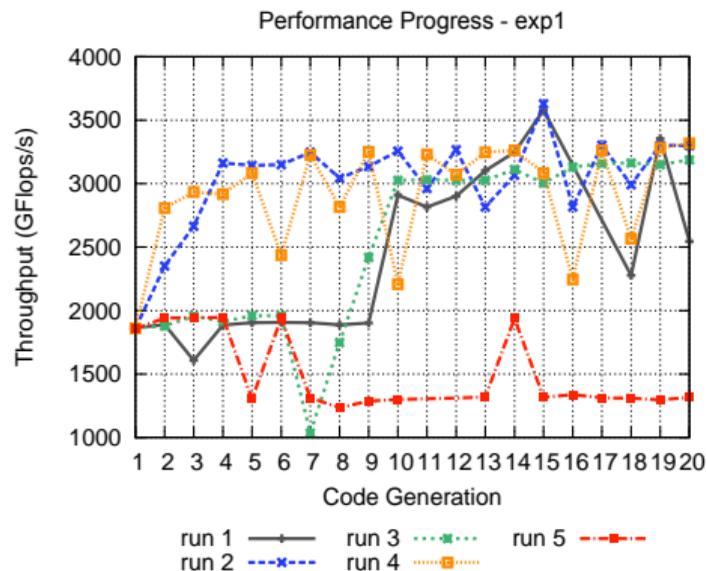


- CSR 形式の疎行列-密行列積 (gridgena.mtx: $n = 48962$, $nnz = 512084$, 右辺幅 $m = 128$) ; 不規則なメモリアクセス ; 疎行列の構造に依存
- 最速 : **78.45 GB/s** (run 5 の 15 世代) (MKL: 86.88 GB/s, **MKL の 90%**)
- FP16 の利用を試行 (失敗), 格納形式の変更を試行 (CSC, ELLPACK, 失敗) → CSR 形式を認識している証拠

DGEMM (CUDA) – 入力コードと実験結果

入力コード

```
1  --global__ void exp1_kernel(int N, const double *
2  A, const double *B, double *C) {
3  int row = blockIdx.y * blockDim.y + threadIdx
4  .y;
5  int col = blockIdx.x * blockDim.x + threadIdx
6  .x;
7  if (row < N && col < N) {
8  double sum = 0.0;
9  for (int k = 0; k < N; k++) {
10     sum += A[row * N + k] * B[k * N + col
11     ];
12     C[row * N + col] = sum;
13 }
14 }
15 extern "C" void exp1_opt_kernel(int N, const
16 double *d_A, const double *d_B, double *d_C) {
17 dim3 blockDim(16, 16);
18 dim3 gridDim((N + blockDim.x - 1) / blockDim.
19 x, (N + blockDim.y - 1) / blockDim.y);
20 exp1_kernel<<<gridDim, blockDim>>>(N, d_A,
21 d_B, d_C);
22 cudaDeviceSynchronize();
23 }
```



NVIDIA Tesla V100

最速コード: 3628.29 GFlops/s

cuBLAS: 6168.09 GFlops/s (cuBLAS
の 59%)

DGEMM (CUDA) – 生成された最速コード (1/3)

```
1 #include <cuda_runtime.h>
2 #include <stdlib.h>
3
4 extern "C" cudaError_t cudaGetDeviceProperties_v2(
5     cudaDeviceProp *prop, int device) {
6     return cudaGetDeviceProperties(prop, device);
7 }
8
9 extern "C" long __isoc23_strtol(const char *nptr, char
10     **endptr, int base) {
11     return strtol(nptr, endptr, base);
12 }
13
14 #define TILE 32 // ×3232 タイル
15 #define K_UNROLL 32
16 // 手動アンロールマクロ
17 #define UNROLL_K(iter)
18     do {
19         double a0 = As[calc][threadIdx.y*2 + 0][iter];
20         double a1 = As[calc][threadIdx.y*2 + 1][iter];
21         double b0 = Bs[calc][iter][threadIdx.x*2 + 0];
22         double b1 = Bs[calc][iter][threadIdx.x*2 + 1];
23         regC[0][0] += a0 * b0;
24         regC[0][1] += a0 * b1;
25         regC[1][0] += a1 * b0;
26         regC[1][1] += a1 * b1;
27     } while (0)
28
29 --global__ void exp1_kernel(int N, const double * __
30     restrict __ A, const double * __restrict __ B
31     , double * __restrict __ C) {
```

```
31     double regC[2][2] = {{0.0, 0.0},
32                         {0.0, 0.0}}; //
33                                     ×22 出力
34
35     --shared__ double As[2][TILE][TILE];
36     --shared__ double Bs[2][TILE][TILE];
37
38     int numTiles = (N + TILE - 1) / TILE;
39
40     int baseRow = blockIdx.y * TILE + threadIdx.y * 2;
41     // 2 行相当
42     int baseCol = blockIdx.x * TILE + threadIdx.x * 2;
43     // 2 列相当
44
45     // ---- 1st タイルを load (ping = 0) ----
46     int ping = 0;
47     for (int dy = 0; dy < 2; ++dy) {
48         int aRow = baseRow + dy;
49         int aCol = 0 * TILE + threadIdx.x * 2;
50         double4 a_pair = ((double4*)(A + aRow * N +
51             aCol));
52         double a0 = a_pair.x;
53         double a1 = a_pair.y;
54         As[ping][threadIdx.y * 2 + dy][threadIdx.x * 2]
55             = a0;
56         As[ping][threadIdx.y * 2 + dy][threadIdx.x * 2 +
57             1] = a1;
58     }
59     for (int dy = 0; dy < 2; ++dy) {
60         int bRow = 0 * TILE + threadIdx.y * 2 + dy;
61         int bCol = baseCol;
62         double4 b_pair = ((double4*)(B + bRow * N +
63             bCol));
64         double b0 = b_pair.x;
65         double b1 = b_pair.y;
66         Bs[ping][threadIdx.y * 2 + dy][threadIdx.x * 2]
67             = b0;
68         Bs[ping][threadIdx.y * 2 + dy][threadIdx.x * 2 +
69             1] = b1;
70     }
71 }
72
73 -- syncthreads();
74
75     // 初期タイルロード
76
77     for (int t = 0; t < numTiles; ++t) {
```

DGEMM (CUDA) – 生成された最速コード (2/3)

```
65     int calc = ping;                                     // 現在計算に使用する
66     バッファ
67     // ---- 計算 (calc バッファ) ----
68     UNROLL_K(0); UNROLL_K(1); UNROLL_K(2);
69     UNROLL_K(3);
70     UNROLL_K(4); UNROLL_K(5); UNROLL_K(6);
71     UNROLL_K(7);
72     UNROLL_K(8); UNROLL_K(9); UNROLL_K(10);
73     UNROLL_K(11);
74     UNROLL_K(12); UNROLL_K(13); UNROLL_K(14);
75     UNROLL_K(15);
76     UNROLL_K(16); UNROLL_K(17); UNROLL_K(18);
77     UNROLL_K(19);
78     UNROLL_K(20); UNROLL_K(21); UNROLL_K(22);
79     UNROLL_K(23);
80     UNROLL_K(24); UNROLL_K(25); UNROLL_K(26);
81     UNROLL_K(27);
82     UNROLL_K(28); UNROLL_K(29); UNROLL_K(30);
83     UNROLL_K(31);
84
85     // 次のタイルを prefetch (ping バッファ) –すでに最後の場合はスキップ
86     ping = 1 - ping;
87
88     // バッファを反転
89     if (t + 1 < numTiles) {
90         for (int dy = 0; dy < 2; ++dy) {
91             int aRow = baseRow + dy;
92             int aCol = (t + 1) * TILE + threadIdx.x
93                 * 2;
94             double4 a_pair = *((double4*)(A + aRow *
95                 N + aCol));
96             double a0 = a_pair.x;
97             double a1 = a_pair.y;
98             As[ping][threadIdx.y * 2 + dy][threadIdx
99                 .x * 2] = a0;
100            As[ping][threadIdx.y * 2 + dy][threadIdx
101                .x * 2 + 1] = a1;
102        }
103        for (int dy = 0; dy < 2; ++dy) {
104            int bRow = (t + 1) * TILE + threadIdx.y
105                * 2 + dy;
106            int bCol = baseCol;
```

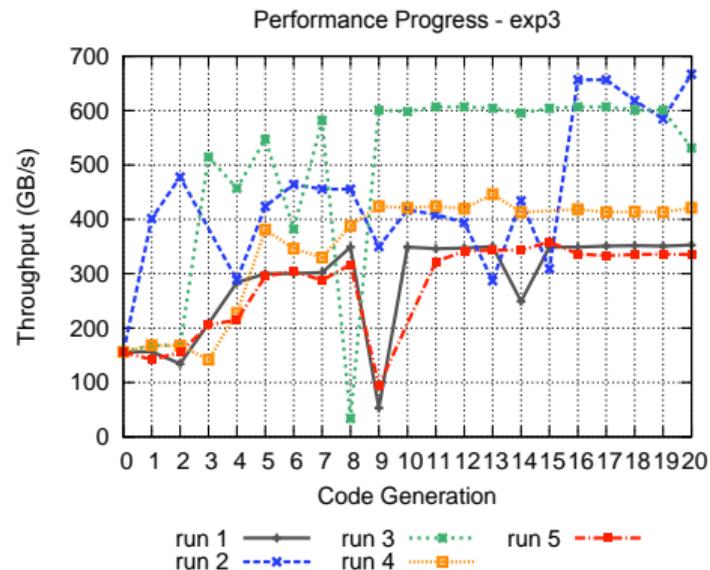
```
92         double4 b_pair = *((double4*)(B + bRow *
93             N + bCol));
94         double b0 = b_pair.x;
95         double b1 = b_pair.y;
96         Bs[ping][threadIdx.y * 2 + dy][threadIdx
97             .x * 2] = b0;
98         Bs[ping][threadIdx.y * 2 + dy][threadIdx
99             .x * 2 + 1] = b1;
100     }
101     }
102     __syncthreads();
103
104     // 次計算に備える
105 }
106
107 // ---- 書き戻し ----
108 C[baseRow * N + baseCol] = regC
109     [0][0];
110 C[baseRow * N + baseCol + 1] = regC
111     [0][1];
112 C[(baseRow + 1) * N + baseCol] = regC
113     [1][0];
114 C[(baseRow + 1) * N + baseCol + 1] = regC
115     [1][1];
116 }
117
118 extern "C" void exp1_opt_kernel(int N, const double *d_A
119     ,
120     const double *d_B,
121     double *d_C) {
122     if (N % 8 != 0) {
123         int Npad = ((N + 7) / 8) * 8;
124         size_t pitch_pad = (size_t)Npad * sizeof(double)
125             ;
126         size_t pitch_orig = (size_t)N * sizeof(double);
127         double *d_Apad = nullptr, *d_Bpad = nullptr, *
128             d_Cpad = nullptr;
129         cudaMalloc(&d_Apad, pitch_pad * Npad);
130         cudaMalloc(&d_Bpad, pitch_pad * Npad);
131         cudaMalloc(&d_Cpad, pitch_pad * Npad);
132         cudaMemset(d_Apad, 0, pitch_pad * Npad);
133         cudaMemset(d_Bpad, 0, pitch_pad * Npad);
134         cudaMemset(d_Cpad, 0, pitch_pad * Npad);
135     }
```

DGEMM (CUDA) – 生成された最速コード (3/3)

```
122     cudaMemcpy2D(d_Apad, pitch_pad, d_A, pitch_orig,  
123                pitch_orig, N,  
                cudaMemcpyDeviceToDevice);  
124     cudaMemcpy2D(d_Bpad, pitch_pad, d_B, pitch_orig,  
125                pitch_orig, N,  
                cudaMemcpyDeviceToDevice);  
126  
127     dim3 blockDim(TILE / 2, TILE / 2);  
128     dim3 gridDim((Npad + TILE - 1) / TILE, (Npad +  
                TILE - 1) / TILE);  
129     exp1_kernel<<<gridDim, blockDim>>>(Npad, d_Apad,  
130                d_Bpad, d_Cpad);  
130     cudaDeviceSynchronize();  
131  
132     cudaMemcpy2D(d_C, pitch_orig, d_Cpad, pitch_pad,  
133                pitch_orig, N,  
                cudaMemcpyDeviceToDevice);  
134  
135     cudaFree(d_Apad);  
136     cudaFree(d_Bpad);  
137     cudaFree(d_Cpad);  
138     return;  
139 }  
140  
141 dim3 blockDim(TILE / 2, TILE / 2);  
142 dim3 gridDim((N + TILE - 1) / TILE, (N + TILE - 1) /  
                TILE);  
143 exp1_kernel<<<gridDim, blockDim>>>(N, d_A, d_B, d_C)  
144 ;  
144     cudaDeviceSynchronize();  
145 }  
146  
147 #undef K_UNROLL  
148 #undef UNROLL_K
```

DCSRMM (CUDA)

```
1  --global__ void exp3_kernel(int m, int n, const
   int *row_ptr, const int *col_indices, const
   double *values, const double *B, double *C) {
2  int row = blockIdx.x * blockDim.x + threadIdx
   .x;
3  int col = blockIdx.y * blockDim.y + threadIdx
   .y;
4  if (row < m && col < n) {
5  double sum = 0.0;
6  int row_start = row_ptr[row];
7  int row_end = row_ptr[row + 1];
8  for (int idx = row_start; idx < row_end;
   idx++) {
9  int sparse_col = col_indices[idx];
10 double val = values[idx];
11 sum += val * B[sparse_col * n + col];
12 }
13 C[row * n + col] = sum;
14 }
15 }
16 extern "C" void exp3_opt_kernel(int m, int k, int
   n, const int *d_row_ptr, const int *
   d_col_indices, const double *d_values, const
   double *d_B, double *d_C) {
17 cudaMemset(d_C, 0, m * n * sizeof(double));
18 dim3 blockDim(16, 16);
19 dim3 gridDim((m + blockDim.x - 1) / blockDim.
   x, (n + blockDim.y - 1) / blockDim.y);
20 exp3_kernel<<<gridDim, blockDim>>>(m, n,
   d_row_ptr, d_col_indices, d_values, d_B,
   d_C);
21 cudaDeviceSynchronize();
22 }
```



Best code: 666.78 GFlops/s

cuSparse: 397.58 GFlops/s

→ 特定の入力に特化した最適化に成功

The AI Scientist (Sakana AI) ^a

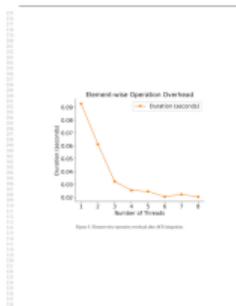
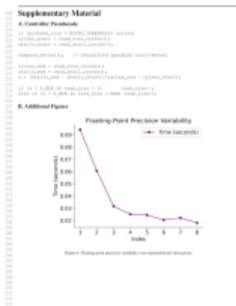
- 研究の自動化を実証 (アイデア生成 → 実装 → 実験 → 評価 → 論文執筆)
- AI 研究に特化, HPC 研究に利用できない

The AI Scientist の HPC 研究者化 ^b

- C/C++実装 → コンパイル → デバッグ・性能評価をするループを実装し, コードの性能最適化を伴う研究に対応
- 自動生成された論文は不十分であるが研究の半自動化を達成しつつある

^aThe AI Scientist-v2: Workshop-Level Automated Scientific Discovery via Agentic Tree Search, arXiv:2504.08066, 2025.

^bTakanori Kotama, Rio Yokota, Daichi Mukunoki, Tetsuya Hoshino, Takahiro Katagiri, Proposal of The AI Scientist v2 for High Performance Computing with Local Large Language Models, SCA/HPCAsia2026, Poster Session, Jan. 2026 (accepted).



商用コーディングエージェント利用の雑感 (1/2)

私の利用状況

- Claude Code, Gemini CLI, Codex の最大プランを契約 (合計月 10 万ぐらい)
- 既存プロジェクト (Fortran/C/C++) の機能追加・最適化・デバッグ・GPU 化で利用
- 最近はほとんどコードを書かなくなった (Vibe Coding)

どういうことができるか？

- 既存コードをほぼ全自動で GPU 移植 (まあまあな性能)
- 専門書に載っているレベルの一般的な数値計算法を理解し最適な手法を提示
 - ▶ コード中のソルバ部分を理解し、疎行列構造を理解し典型的アルゴリズム, 前処理, 疎行列フォーマットを提案してくれる
- ライブラリの利用方法を知っていて自動で環境構築
- コードの妥当性検証のための簡単なテストの生成

商用コーディングエージェント利用の雑感 (2/2)

課題

- ぐちゃぐちゃで大規模なコードほど難易度が上がる
- 嘘, 失敗, 記憶喪失, 堂々巡り, etc.
- 著作権, 責任の所在, 過度な AI 依存, etc.

“Vibe Resarch” 的な利用

- コーディングだけでなくアイデア出しや論文執筆への利用を試行
- 半自動の The AI Scientist for HPC
- 人間が大まかなアイデアを提示 → (ここから LLM エージェント) 背景の妥当性, 関連研究調査, 新規性の検討 → 実装・実験方法を含めて論文を書かせる → それを仕様書としてプログラムを実装し実験 → 論文に反映

素人が使うツールではなく専門家が使う武器

- 全然使えないと言ってる人もいる (特に達人級のプログラマ) が, 使い方次第
- リテラシとして使いこなしを考える
- 急速な技術進化についていく

HPC-GENIE プロジェクト（名古屋大学情報基盤センター）

- 生成 AI の HPC コード開発への活用に向けた研究
- 特にエージェント開発にフォーカス – VibeCodeHPC
- 国産・オープン・ローカル実行可能な技術開発

HPC コード開発における生成 AI への期待

- 課題や制約はまだ多くあるが、研究開発の効率を劇的に向上させている
- **富岳 NEXT を見据えたコードの GPU 化において強い期待感**

提言（議題？）

- **資金** – AI の利用のため、人材の育成のため、人材の確保のため
- 急速な技術発展・状況変化にどうついていくか
 - ▶ 数ヶ月で技術が 1 周している印象
 - ▶ アカデミアにおける教育・研究のやり方自体がついていけない