

# アドバンス・シミュレーション 第6回 ・セミナー 2024

2024年 9月20日(金) 離

## プログラム

1. アドバンスソフト株式会社についてと、先生のご紹介 ..... 1

### 招待講演

2. 「ポスト・エクサ、ポストムーア時代の HPC と AI」 ..... 3

神戸大学大学院 理学研究科 惑星学専攻 特命教授 牧野 淳一郎 様

3. アドバンスソフトからの情報提供 ..... 43

## 講演概要

### 「ポスト・エクサ、ポストムーア時代の HPC と AI」

神戸大学大学院 理学研究科 惑星学専攻 特命教授 牧野 淳一郎 様

「富岳」を始めとするエクサスケールシステムは世界的にも立ち上げが済み、実際の運用に入っている。また、日本ではポスト富岳のフィージビリティスタディが進み、次世代計算基盤の姿が見え始めてきた。一方、人工知能、特に生成 AI・大規模言語モデルの応用が急激に普及し、汎用 HPC、さらには汎用 CPU の市場よりも AI 向け GPGPU の市場のほうが大きくなっており、HPC 向けのプロセッサの将来像は見えにくくなっている。また、半導体技術は完全にポストムーア時代に入り、開発費は高騰する一方、世代毎の性能向上はマージナルになっている。

この状況下でのプロセッサと HPC アプリケーションの将来像を展望する。



アドバンスソフト株式会社 セミナー事務局

〒101-0062 東京都千代田区神田駿河台四丁目3番地 新お茶の水ビルディング 17階西

TEL: 03-6826-3971 FAX: 03-5283-6580

URL: <http://www.advancesoft.jp/> E-mail: [office@advancesoft.jp](mailto:office@advancesoft.jp)

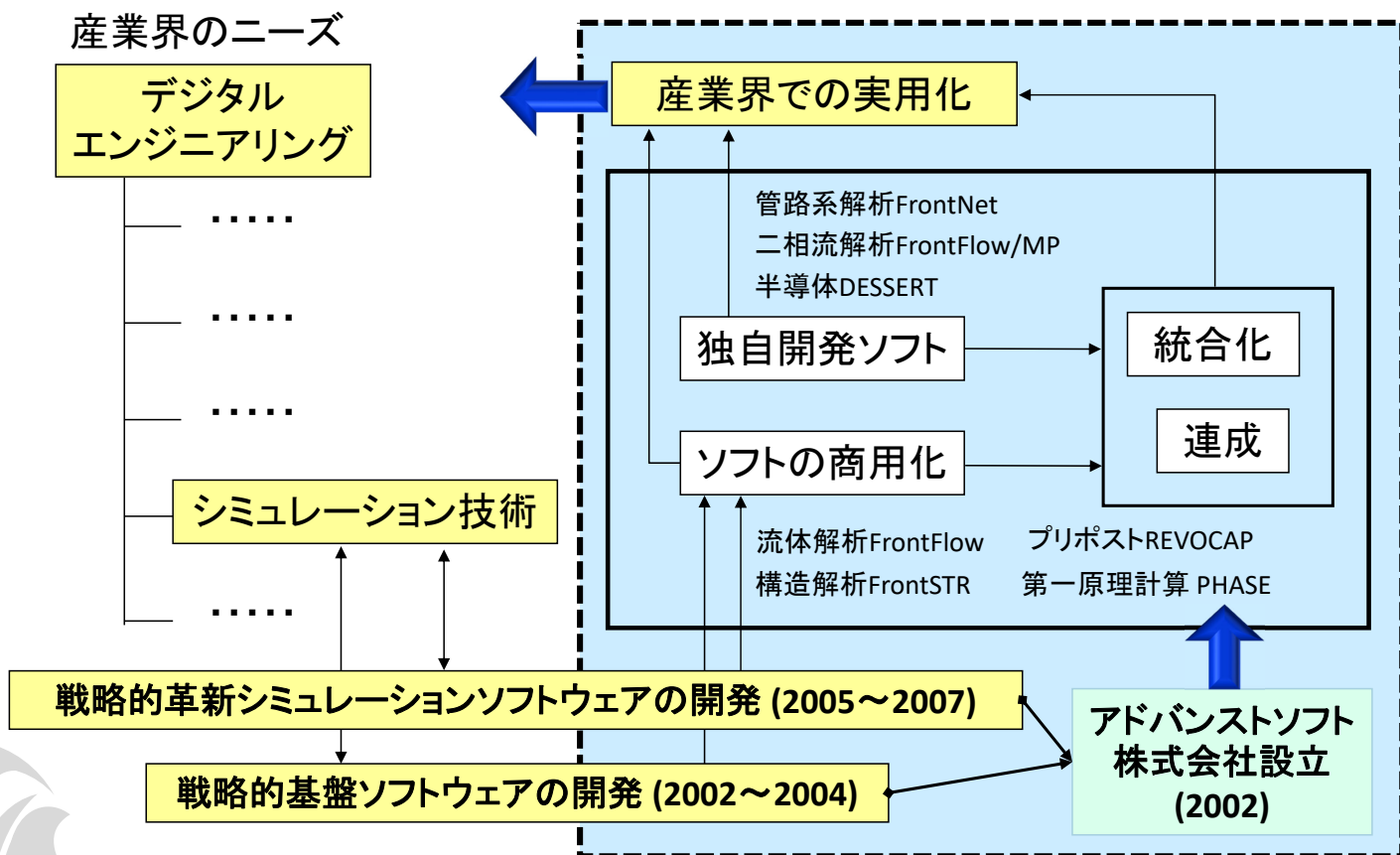


# アドバンスソフト株式会社についてと 先生のご紹介

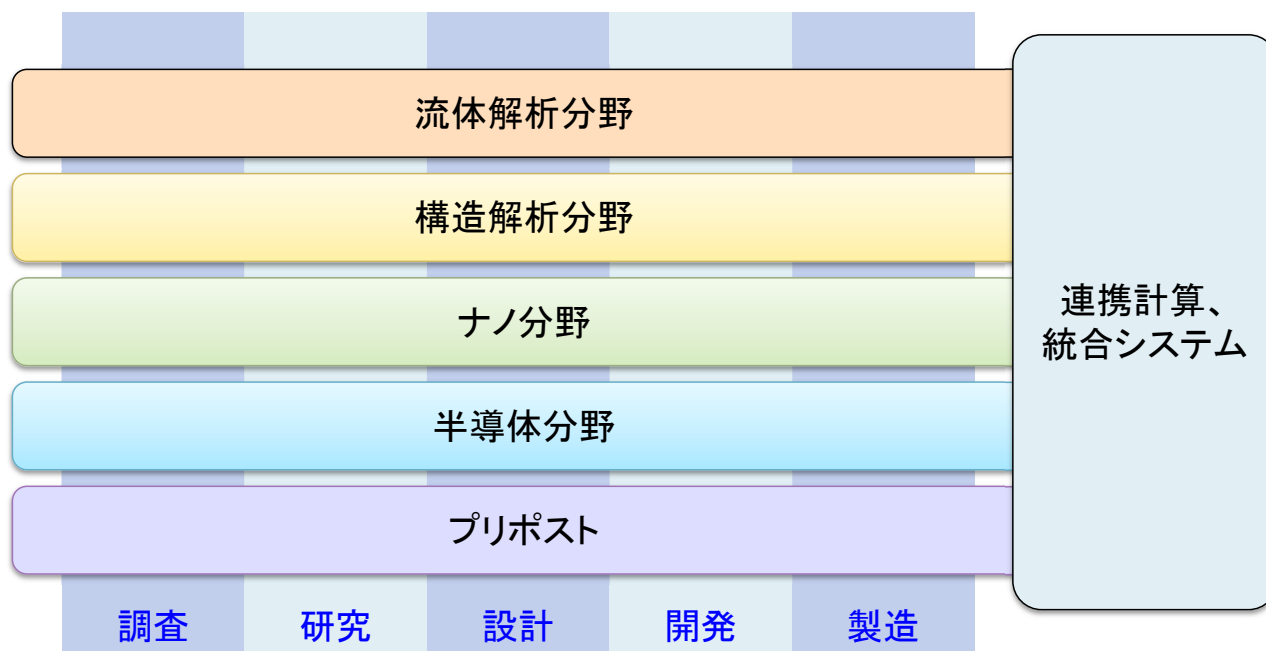
第6回 アドバンス・シミュレーション・セミナー 2024  
(神戸大学大学院 理学研究科 惑星学専攻  
特命教授 牧野 淳一郎様)

2024年9月20日 (金) 開催  
アドバンスソフト株式会社

## アドバンスソフトとは



# 事業分野



産業の主要な分野のあらゆるフェーズで直面する課題に対し、  
科学技術計算によるソリューションをご提供します。

# 牧野 淳一郎 先生 のご紹介

## ご経歴・ご研究内容

1990年	東京大学大学院 総合文化研究科 広域科学専攻博士課程 修了（学術博士）。
1999～2006年	東京大学 助教授。
2006～2011年	国立天文台 教授。
2011～2013年	東京工業大学 教授。
2014～2021年	理化学研究所 計算科学研究機構 副プロジェクトリーダー。
2016～2023年	神戸大学 教授 等を経て、
2023年～	神戸大学 特命教授 と、株式会社 Preferred Networks VP コンピュータアーキテクチャ担当 CTOを兼任。

# ポスト・エクサ、ポストムーア時代の HPC と AI

牧野淳一郎

神戸大学理学研究科惑星学専攻

第6回 アドバンス・シミュレーション・セミナー

2024/9/20

## 話の構成

- これまでの計算機の進歩と現状の困難
- 「ポストムーア」
- 計算機の進歩の方向
- AI 向けプロセッサの方向
- **MN-Core** について
- シミュレーション用計算機は？
- アプリケーションは？
- まとめ

# 牧野はどういう人か？

- 元々、銀河系や星団の、コンピュータシミュレーションを使った研究をしていた(今もしている)
- 「シミュレーション専用計算機」の開発を 1990 年くらいから始める
- 2004 年から、汎用並列計算機開発にもかかわる
- 2011 年から、理研の「ポスト京」プロジェクトにも
- 2016 年から、AI ベンチャーの Preferred Networks と共同で深層学習向けプロセッサを開発
- 2020 年 6 月、完成したプロセッサ MN-Core を使った MN-3 が「Green500」ランキングで世界一位に
- 2023 年 11 月から Preferred Networks VP 計算アーキテクチャ担当 CTO 兼任(クロスアポイントメント)

## これまでの計算機の進歩

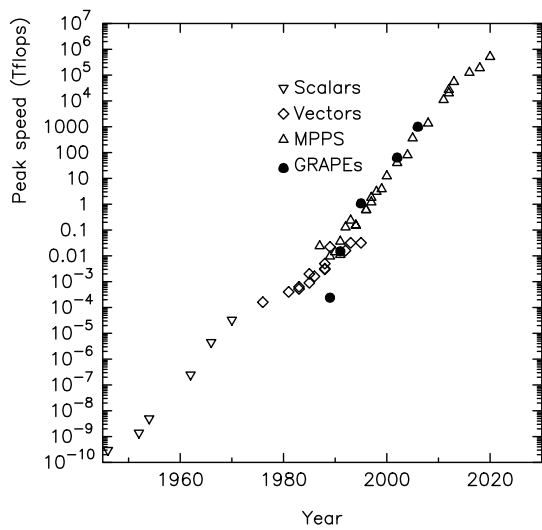
1940年代から2010年代までの70年間、ほぼ10年で100倍。何故そのような指数関数的進歩を長期に続けたのか？

基本的な理由:

- 使うスイッチ素子が高速になった
- 使うスイッチ素子が小型、低消費電力になって、沢山使えるようになった
- 使うスイッチ素子が安くなって、沢山使えるようになった(スパコンの物理的大きさは70年代が最小。そこまで段々小さくなって、そこからまた大きくなった)



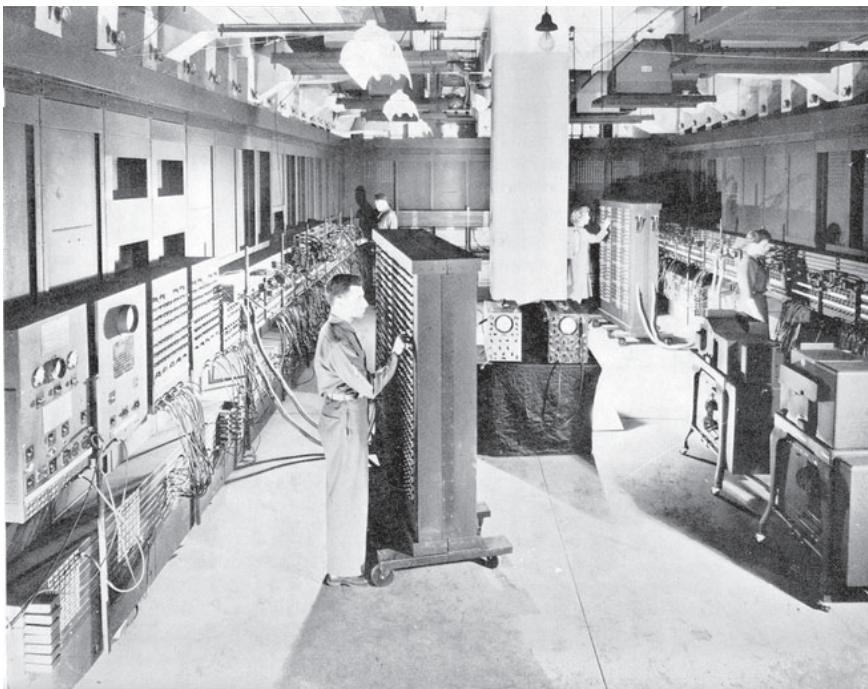
# スパコンとは何か



- 単純には、「その時代で最高速クラスの計算機」
- 80年間で16桁弱速くなった=ほぼ10年で100倍
- 何故これほど進歩したか?は「スパコンとは何か」そのもの

(但し、現状では **Top500** の数字があんまり意味がなくなっているかも、、、 **GPT** とかの **LLM** 学習に使われているシステムのほうが大きい)

## 1940年代の「スパコン」



**ENIAC**

# ENIAC

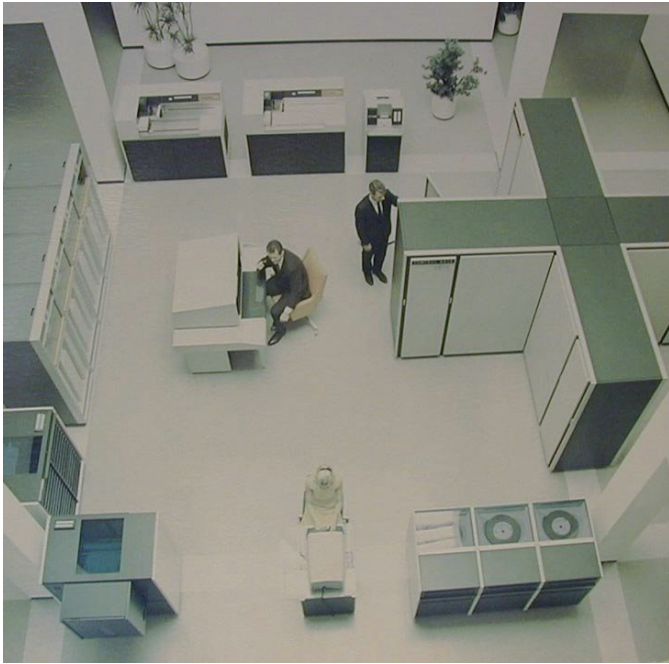
- 第二次世界大戦中に、アメリカ陸軍が弾道計算のために開発。完成は**1946年**
- **18,000**本の真空管、消費電力**150KW**、**10**進法**10**桁の数が基本
- 「プログラム」はスイッチ、ケーブルの変更で行う
- 加減算が**1秒5000**回、乗算が**400**回くらい(らしい)

## 1950年代の「スパコン」



IBM 709 — IBM 最後の真空管計算機

## 1960年代の「スパコン」



**CDC 6600 1Mflops**

### **CDC6600**

- 1964 年から出荷。シリコントランジスタを利用
- 10MHz クロック、1Mflops 程度の性能
- 「スカラー計算機」1つの命令は演算やデータ移動1つを記述
- 近代的計算機アーキテクチャの原型
  - ロード・ストアアーキテクチャ
  - 複数の演算器の「アウト・オブ・オーダー」実行
- シーモア・クレイが設計

## 1970年代の「スパコン」



**CRI Cray-1 160Mflops**

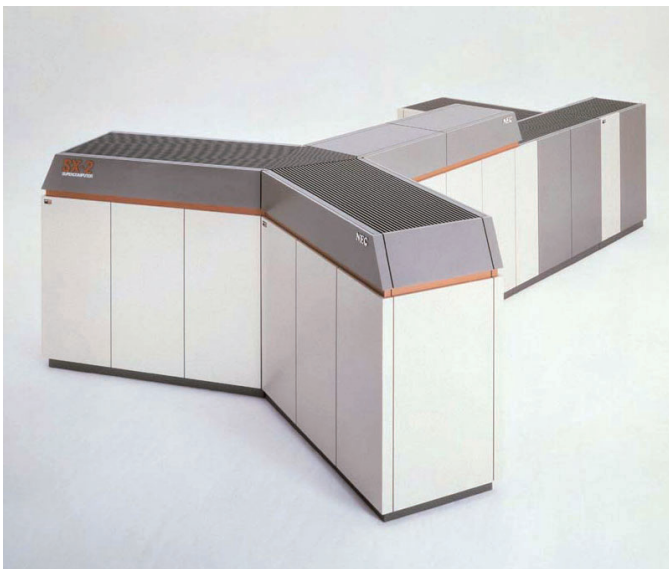
### **Cray-1**

- 1976 年から出荷。IC (ECL 素子による小規模 IC) を使用
- 80MHz クロック、160Mflops
- 最初に成功したベクトル計算機
  - ベクトルレジスタ
  - 半導体メモリ
- シーモア・クレイが設計

## ベクトル計算機って？

- 「スカラー計算機」と「ベクトル計算機」と分ける
- 「ベクトル処理」をするかどうか
- スカラー計算機(あるいは普通の計算機の基本命令): 命令1つで演算1つ。例えば掛け算1つとか
- ベクトル計算機: 1命令で複数の演算を処理。ベクトルの、要素同士の四則演算が基本。それだけでは足りないので色々追加機能がある
- 必ずしも並列処理するわけではない。Cray-1では「パイプライン処理」。長さ  $n$  のベクトル命令の実行に 定数 +  $n$ クロックかかる

## 1980年代の「スパコン」

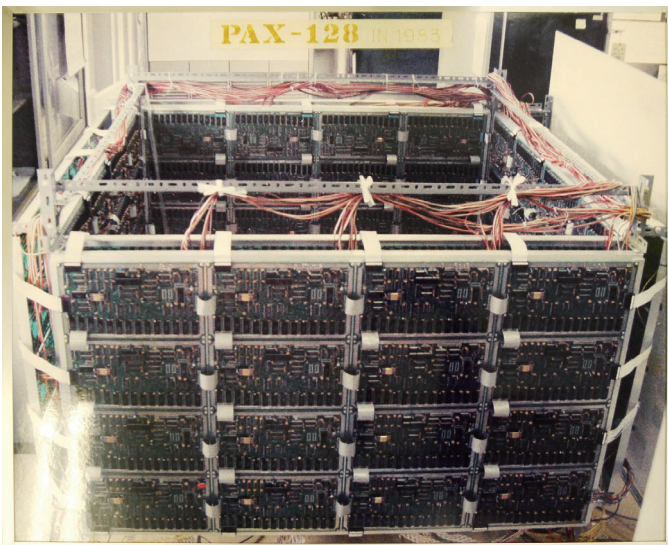


NEC SX-2 1.3Gflops

## NEC SX-2

- 1985 年から出荷。CML 論理素子による LSI を使用。1000 ゲート程度 (Cray-1 の IC の 100 倍程度)
- 6ns クロック、掛け算、加減算パイプラインをそれぞれ 4 本
- 富士通、日立が同様なマシンを出荷、NEC は最後
- Cray は共有メモリマルチプロセッサに (Cray XMP, YMP)

## 1980年代にはこんなのも



PAX-128 4Mflops。筑波大学、星野ら

## 1990年代の「スパコン」



**ASCI-Red, Intel Pentium Pro, 1.8TF**

## ASCI-Red

- 1996年に完成。200MHz Pentium Pro プロセッサ 9216台
- 3次元メッシュネットワーク(ほぼ2次元、、、)
- 核爆弾のシミュレーション用

## 2000年代の「スパコン」



**IBM BG/L 360TF**

### **IBM BG/L**

- **2004**年完成、カスタムプロセッサを最大 **131072**個使用。**3次元**トーラスネットワーク
- ピーク性能**360TF** (もっと大きい構成もあった模様)
- 後継の **BG/P**, アーキテクチャを一新した **BG/Q** の開発後、プロジェクト解散

## 2010年代の「スパコン」



神戸市・ポートアイランドの「京」コンピュータ、10PF

### 「京」

- 2012年完成、カスタムプロセッサを8万個使用。6次元トーラスネットワーク
- ピーク性能11PF
- 商用版富士通FX10, 「FX100」が開発された。2020年「富岳」への入れ換えが完了
- 「富岳」は「京」の約50倍のスピード
- (牧野は開発プロジェクト副プロジェクトリーダーの1人)

## 現在(2020年以降)の普通のスパコン

- NVIDIA の GPU を使用
- 沢山並べる
- 高速ネットワークを使う、あるいは Cray 社のシステムを買う
- 独自プロセッサとかはあまり使われない
- あんまり芸がないので写真は省略

## 過去のスパコンの進化

何の話をしたかったかというと：何故計算機はどんどん速くなったのか？

基本的な理由：

- 使うスイッチ素子が高速になった
- 使うスイッチ素子が小型、低消費電力になって、沢山使えるようになった
- 使うスイッチ素子が安くなって、沢山使えるようになった(スパコンの物理的大きさは70年代が最小。そこまで段々小さくなって、そこからまた大きくなった)

## 素子の高速化

- といっても、真空管でもそれなりに速かった
- スイッチング速度が重要でないわけではないが、配線を信号が伝搬する速度のほうが昔から重要
- 昔は信号はほぼ光の速さでつたわった
- 最近の LSI 上の配線は非常に細く (抵抗が大きく)、キャパシタンスを充電しないといけないことによる RC(抵抗とキャパシタンス) 遅延のため、信号が伝わる速度は光速度よりはるかに低い
- 太い配線に大電流を流せば速いが、大量の電力消費になる

## 素子の小型化

- 真空管 → トランジスタ → IC という進化は 1970 年代までは重要
- サイズだけでなく、消費電力が下がることが重要
- 80 年代から重要になったのは CMOS LSI の微細化。10 年でサイズが 1/10 になる
- CMOS 素子では(2000 年くらいまでは)微細化すると電圧を下げる  
ことができ、消費電力が下がり、速度は向上した。いわゆる CMOS  
スケーリング。
- 2000 年頃からは電圧が下がらないので、電力はちょっと下がるが  
速度は上がらなくなった。いわゆる CMOS スケーリングの終焉。
- そろそろ微細化も困難になってきた。また、トランジスタの構造・  
製造工程が複雑になり、微細化するとかえって価格上昇するよう  
になった。いわゆるムーアの法則の終焉。

# CMOS スケーリングって何？

「Dennard Scaling」とも

トランジスタのサイズ (3 次的にすべて)、電源電圧を  $1/k$  にし、不純物濃度を  $k$  倍にすると、トランジスタの速度は  $k$  倍、スイッチングあたりの消費電力は  $1/k^3$  になる

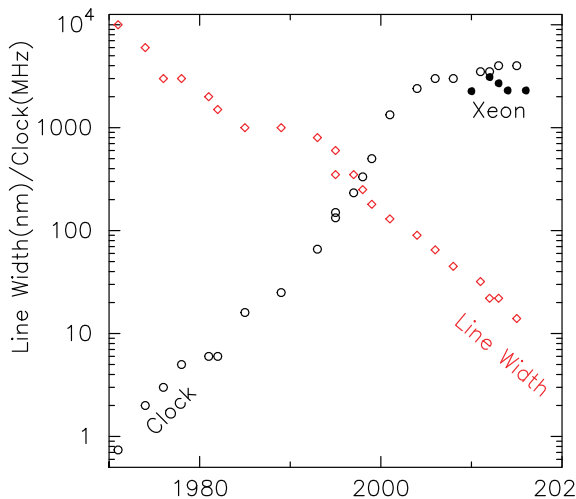
これが成り立つなら、

動作周波数  $\propto 1/\text{線幅}$

動作電圧  $\propto \text{線幅}$

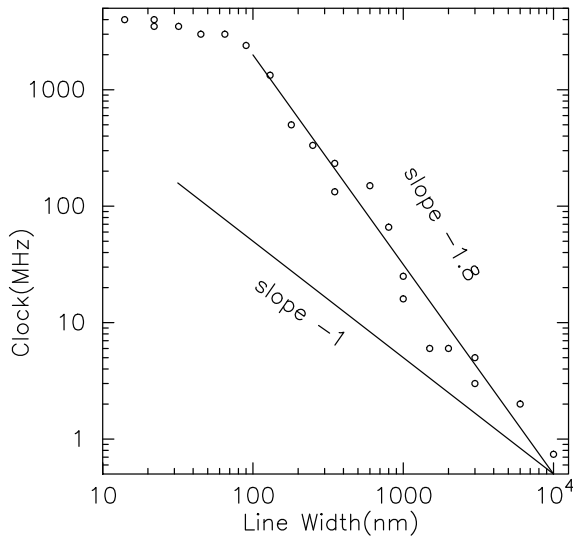
性能  $\propto 1/\text{線幅}^3$  (電力、面積一定で)

## 現実には



40年間の Intel マイクロプロセッサの線幅(トランジスタサイズ)と動作クロック  
線幅は40年間で3桁縮小。クロックは25年間で3桁上昇、そのあと停止。(コア数の多い Xeon は段々クロック下がる)

## 線幅とクロック



線幅とクロックの関係

**90nm** までは

クロック  $\propto$  線幅<sup>-1.8</sup>

そのあとはほぼクロック一定

(Xeon は下がる、、、)

**90nm** までのクロック上昇は理屈より速い=無理していた。

## つまり

- ムーアの法則(トランジスタサイズは時間の指数関数)は**2015**年くらいまでは成り立っていた(今は成り立たってない)
- いわゆる **CMOS** スケーリング(速度が線幅に反比例)は**2000**年くらいまでしか成り立たってない

何が起こったのか？

- **90nm** までは電圧をあまり下げず、パイプライン段数を増やすことでクロックをあげた
- これは消費電力増やすので、**1チップ100W**になったところで限界
- そこからは、コア数やコア内演算器数を増やすことで性能向上
- 最近はそれも難しくなってきた再びクロック上げる方向も

ちょっといきあたりばったり観あり。というわけで、、、

## 計算機アーキテクチャの進化を振り返ってみる

- 1976 年まで: スカラー計算機。最後は CDC 7600
- 1976 年から 1992 年まで: ベクトル(共有メモリ並列含む)計算機。Cray-1 から C-90 まで
- 1993 年から 2008 年まで: (マルチコア含む) マイクロプロセッサの分散メモリ並列計算機 Cray T3D から、、、 Red Storm あたりまで
- 2008 年から: GPU アーキテクチャのアクセラレータとマイクロプロセッサの組合せ: IBM RoadRunner (Cell なので GPU じゃないけど)

大体 15 年毎にアーキテクチャが大きく変わった。GPU の次はまだでてきてない

## アーキテクチャの変化が必要な理由

基本的な理由: そのアーキテクチャでは、半導体の進歩を有効に利用できなくなる

- 半導体技術の変化
- アーキテクチャのスケラビリティの限界

## スカラー計算機 ベクトル計算機

- スカラー計算機の進歩: 増えるトランジスタを「1つの演算器を速くする」に
- 主記憶は磁気コアで、半導体よりはるかに遅い
- **CDC 7600 (S. Cray の設計)** あたりが最後
- **S. Cray** はこの後、4プロセッサ並列の **CDC 8600** を開発していたが、完成前に **CDC** やめて **Cray Research** を設立。ベクトルの **Cray-1** を開発する

スカラー計算機では

- **IC** の開発で使えるゲート数が1演算器を超えて大きくなった
- 半導体メモリの開発でメモリが非常に高速になった

ことを生かせなかった

ベクトル: 半導体メモリは有効利用できた

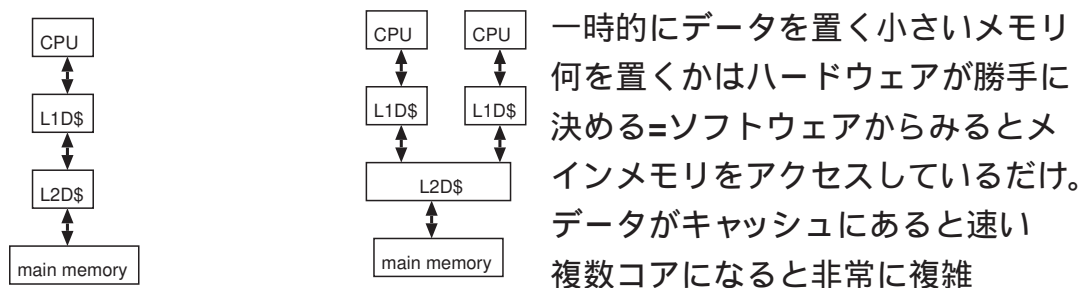
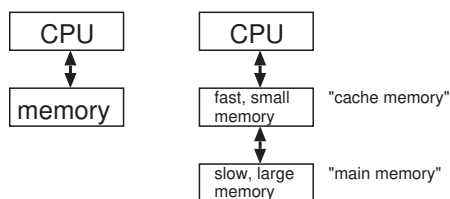
## ベクトル計算機 分散メモリ並列計算機

- ベクトル計算機の進歩: 1プロセッサのパイプライン数や主記憶を共有するプロセッサの数を増やす
- パイプライン数、プロセッサ数に比例以上にメモリへの配線の数が増える。ある程度 (**64**パイプラインあたり) で限界
- 多数の小さなプロセッサとメモリの組合せの間を細いネットワークでつなぐ構成が有利になる
- 初期の代表: **Cray T3D**。但しネットワークが太い設計でそのうち破綻。

# マイクロプロセッサ超並列 GPU

- 1チップの中に沢山プロセッサがはいるようになる
- これは実はシステムレベルでのベクトル計算機の限界と同じ
- マルチ(メニー)コアプロセッサでは、キャッシュコヒーレンシを維持した階層キャッシュとオフチップの共有メモリ
- キャッシュコヒーレンシの維持のための電力が支配的になる
- GPU ではキャッシュコヒーレンシを緩和したり、捨てたりで若干改善
- 本質的な解決ではない

## キャッシュとは？

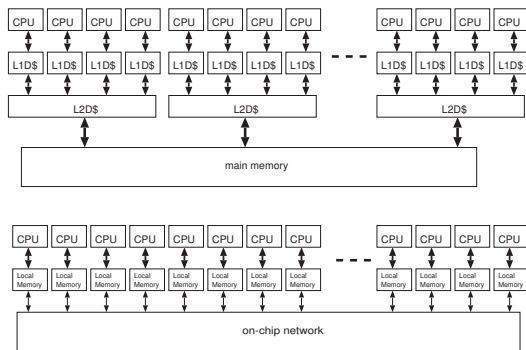


一時的にデータを置く小さいメモリ  
何を置くかはハードウェアが勝手に決める=ソフトウェアからみるとメインメモリをアクセスしているだけ。  
データがキャッシュにあると速い  
複数コアになると非常に複雑

あるCPUが自分のキャッシュに書いたものが他のCPUが元々のそのデータのアドレスを読んだ時にちゃんと読める必要がある(キャッシュコヒーレンシ)

# GPU ???

- コヒーレンスを諦めても、オフチップメモリや階層キャッシュで物理的に遠くにデータ移動させること自体が性能の制限になる
- と書いた以上、階層キャッシュを諦めて、なるべく物理的に遠くにデータ移動させないようにすることが必要になるはず
- つまり、1チップの中でも、多数の独立なプロセッサに分かれた、チップ内分散メモリに行く必要がある



## 何故横方向のデータ移動が制約になるか？

### 配線遅延・配線消費電力の問題

微細化しても、「単位長さ辺りのキャパシタンス」はかわらない、ところが、「単位長さあたりの抵抗」は配線幅の二乗に反比例して増える:

- 配線長がスケールしても配線遅延は「デザインルールに反比例して」増える
- 配線長一定だと遅延は二乗で増える。電力は減らない

(最近だと、トランジスタ密度は上がっても配線はそんなに細くならないので配線が短くなれば遅延・電力が減る。長いと電力も遅延も減らない)

# 数字をいれてみると

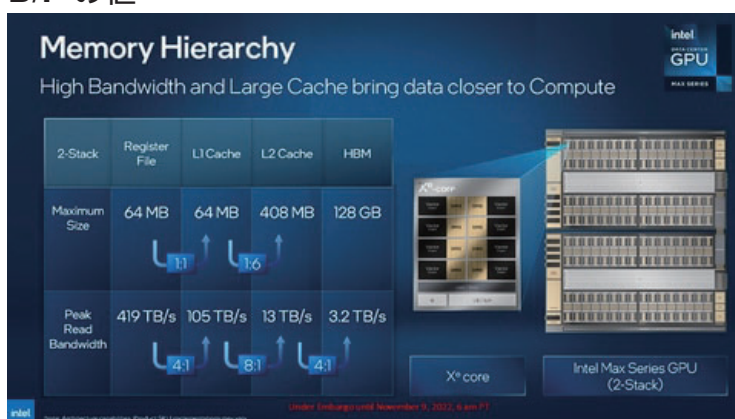
- 10cm ( $10^5 \mu\text{m}$ ) の配線は 20pF のキャパシタンスがある。これは微細化してもかわらない(線と平面の間のキャパシタンス)。電圧 1V だと、この線は 10pJ/bit を消費する
- DDR5 メモリは 20pJ/bit くらい。電圧 1.3V とか。
- LPDDR5 と GDDR6x: 10 pJ/bit くらい。配線が DDR5 のモジュールより短くて電圧も低い
- HBMx: 3-4 pJ/bit。概ね 25mm くらいまで配線短くなっている

最先端の GPU だと、メモリアクセスだけで消費電力の半分くらいになる

## Intel GPU の例

B/F: 1 演算の間に何バイトのデータを移動できるか

B/F の値



レジスタ: 8? (普通 16 いる。行列乗算器だと 8 はありえる)

L1: 2

L2: 0.25

メモリ 0.06

非常に小さい

それでも電力消費が大きい原因になっている

# 最新の GPU の状況

	FP64 性能 (TF)	L2 バンド幅 (B/F)	メモリバンド幅 (B/F)
V100	7.8	0.32	0.12
A100	19.5	0.25	0.08
AMD MI250X	47.9	0.27	0.07
H100	34	0.17	0.10
Intel MAX	52(31)	0.25	0.06

- L2 バンド幅が極端に小さくなっている。0.2 前後だと L2 にのるアプリケーションでもメモリアクセスリミットなら実行効率は 5% 程度。
- CPU に比べてメモリバンド幅・演算性能の絶対値は高いが、キャッシュでも B/F は低く、アプリケーション実行効率が高くない。

## ポスト GPU アーキテクチャ

ポスト GPU アーキテクチャでは

- チップ内部の演算器が物理的に共有するオフチップメモリはもたない
- チップ内部の演算器が物理的に共有するキャッシュはもたない

必要がある。これは 1990 年代前半の共有メモリベクトル並列から分散メモリ超並列への移行と同じ。

つまり、チップ内で分散メモリ超並列を実現すればよい。

但し、普通にオフチップ DRAM をつけたのでは駄目 (バンド幅が足りない)

普通でないオフチップ DRAM= 3D 実装

この話の前に深層学習向けアクセラレータの話

## 深層学習向けの計算機

- ニューラルネットワークの1層は「行列ベクトル積」
- 画像認識で主流の「畳み込みネットワーク」では「行列行列積」
- 計算精度はあんまりいらぬ。最近は8ビット表現、さらに4ビット表現も使われる
- トランスフォーマーでも「行列行列積」が主要な計算

## 現在の最先端

- **NVIDIA H100 (2022年発表)**
- **16ビット**で  $4 \times 4$  の行列同士の乗算の専用回路、**8ビット**だとその**2倍**の性能になるプロセッサ
- **1チップで 1 Pflops** 程度を実現。(「富岳」の1チップは **11 Tflops** なので、**100倍**近く速い)
- 通常の**64ビット**だと **50 Tflops** 程度。
- 今年 **B200** がでる予定。性能**5倍**だがプロセッサダイが**2個**はいついて電力も**2倍**近い。

以下、Preferred Networks でやっている話を

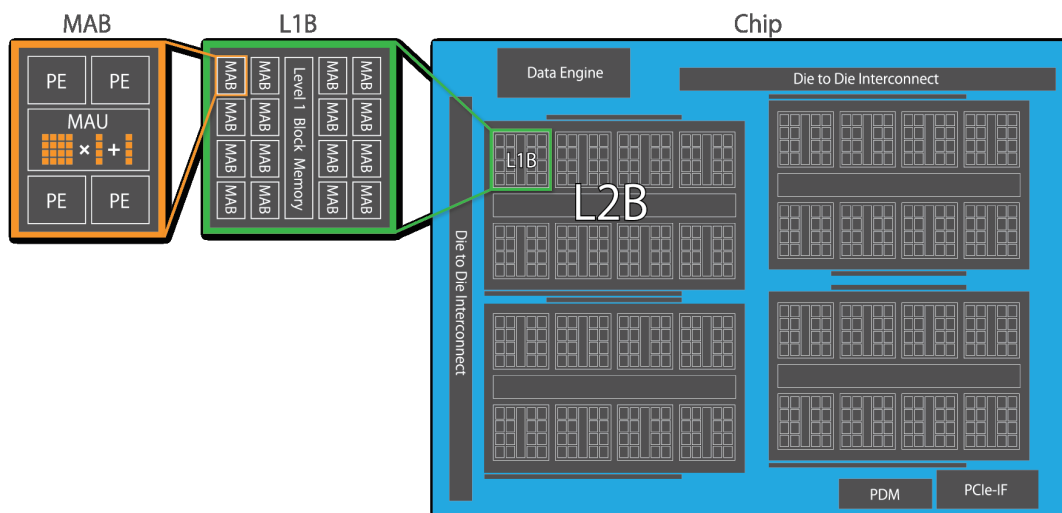
# MN-Core

- **PFN** と牧野のところで開発した深層学習向けプロセッサ
- 深層ニューラルネットワーク (**DNN**) の学習で世界最高の電力あたり性能を実現
- **2020** 年に完成
- **FP16** (ライク) フォーマットでピーク **524TF**
- 電力消費 (チップレベル) **500W** 以下、**1.2TF/W**
- 同じ半導体技術の **NVIDIA V100** の **2.5** 倍以上の電力あたり性能を実現した

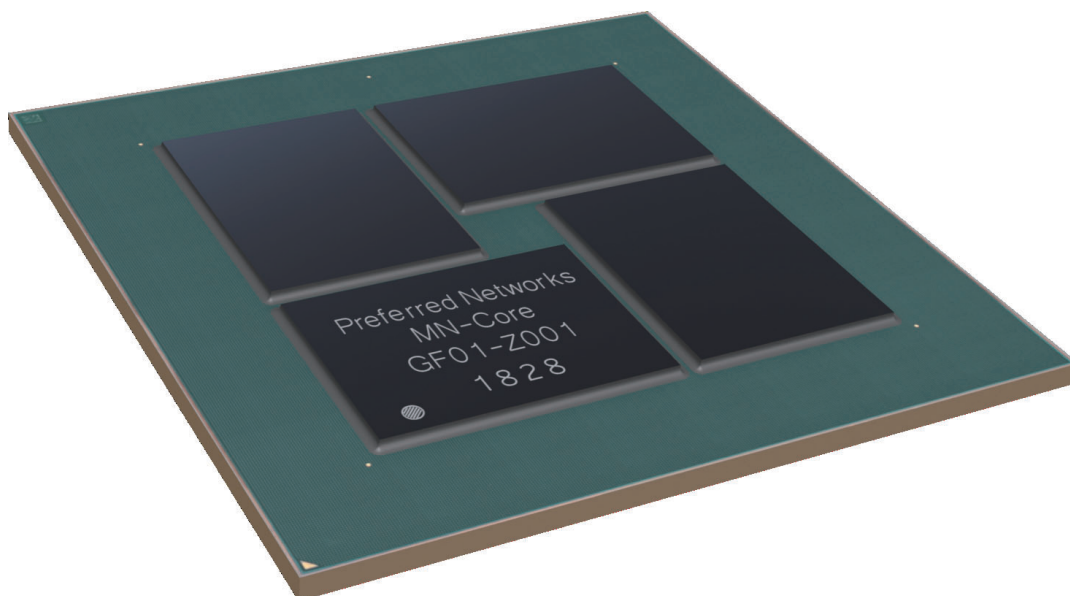
## MN-Core 概要

- 1 カード: 1 「モジュール」
- 1 モジュール: 4 ダイ **MCM**
- 1 チップ: **PCIe (gen4, x16)**, **xDDRx** メモリ (内緒) 4 “**Level-2 放送ブロック**” (**L2Bs**)
- 1 **L2B**: 8 **L1Bs**
- 1 **L1B**: 16 **MABs (Matrix Arithmetic Blocks)**
- 1 **MAB**: 4 **Processor Elements**。これに対して 1 つ 「行列乗算ユニット」
- **FP64:FP32:FP16** 性能比は **1:4:16**
- カード全体が **SIMD** で動く。命令ストリームは 1 つ。

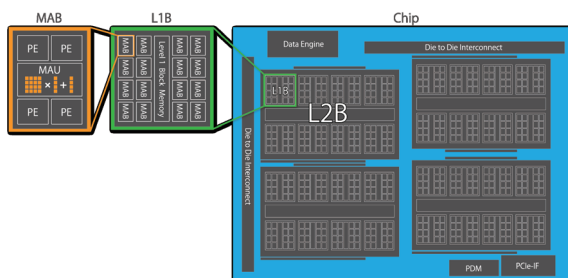
# MN-Core 構成図



# MN-Core



# チップ階層構造



- トップレベル: 4チップ内に4個(合計16個)のL2B(レベル2放送ブロック)
- L2Bの中: 8個のL1B(レベル1放送ブロック)
- L1Bの中: 16個のMAB(行列演算ブロック)

- MABの中: 1つのMAU(行列演算ユニット)と4個のPE(プロセッシングエレメント)

**MN-Core:** 1ボードにPE 8192個。MAB 2048個。ポスト富岳の現在の案はPEあたりの性能を2倍に強化、クロック5倍、PE数1.5倍くらい。

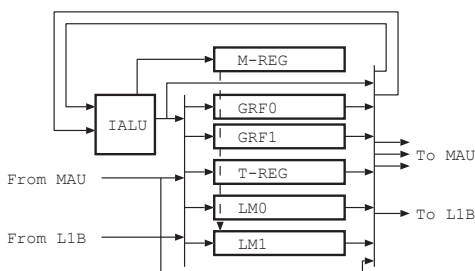
## チップ階層構造続き

- 概念的には全体が同期して単一の命令流を処理する。L2B以下は実際にクロック同期。チップ間を含めてL2B間は(周波数は同じだが)命令実行にずれがあることを許す。
- チップ毎にDRAMブロックがある。DRAMはL2Bとブロック転送でデータ交換できる。放送、個別転送、チップ内分配等いくつかのモードがある。

# 詳細構造

- PE(プロセッシングエレメント)
- L1B(レベル1放送ブロック)
- L2B(レベル2放送ブロック)
- 全体

## PE(プロセッシングエレメント)



- 演算器は IALU と MAU
- MAU は倍精度、単精度、半精度精度の行列ベクトル積を実行
- L2B の中: 8 個の L1B(レベル1放送ブロック)
- L1B の中: 16 個の MAB(行列演算ブロック)

- バスみたいに書いてますが実際の回路はマルチプレクサ
- GRF は (MN-Core では) 1R1W 2ポート。LMx はシングルポート。
- T-reg: 補助レジスタ。1R1W で1ベクトル分
- LMx (ローカルメモリ) は 2048 64 ビット語 x 2、GRF は 512 語
- ベクトル長 4 の固定長ベクトル命令。

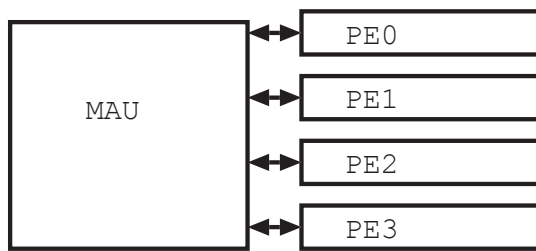
## PE(プロセッシングエレメント) 続き

- 直前の命令の演算結果を入力オペランドにとることが可能(フィードバック)
- ベクトル命令でのアドレッシングは固定、連続アドレス、ストライドアクセスが可能(ベクトルレジスタもストライドアクセス可能)
- T-レジスタを使った間接アドレスも可能
- アドレスはベースアドレスレジスタで修飾
- IALU の演算結果から大小比較、一致比較等のフラグベクトルを生成して M レジスタに格納可能
- M レジスタの値でメモリ、レジスタに書き込むかどうかを制御

## PE(プロセッシングエレメント) 続き

- ロードストアアーキテクチャではなく、LMx の読み出し結果を直接演算器の入力にできる。また、演算器の出力をメモリに直接格納できる。複数のユニットに格納もできる。
- B/F が極端に大きい。(ベクトル演算モードに対して 8 ある)。このため、ローカルメモリにデータがあればほぼ理論ピーク性能がでる。
- この機能と、演算結果フィードバック、Tレジスタを最大限利用することで、通常のアーキテクチャに比べてレジスタファイルのハードウェア規模、消費電力を大きく削減している。また、スーパースカラー実行制御、OoO 実行のための回路も不要となる。=高い電力性能の実現に大きく貢献

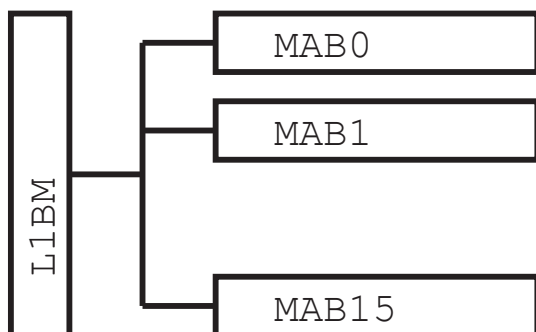
## MAB(行列演算ブロック)



- PE4 個が1つの行列乗算器を共有
- サイクル毎に行列ベクトル積を実行。  
 $d = A * b + c$
- MAU は行列レジスタをもつ

- 行列レジスタは複数あり、演算中に別の面への書き込みが可能
- 転置しながら格納も可能
- 倍精度、単精度、半精度行列に対応。サイクルあたり演算数は **32, 128, 512**
- ベクトル演算 (PE から見ると **64** ビットスカラー演算) モードもある
- 倍精度、単精度、半精度行列乗算器は同一の乗算器の構成変更で実現 (特許取得済): 他社の回路に比べると同一性能で回路規模が半分以下

## L1B(レベル1放送ブロック)



- MAB 16 個が1つの L1BM(レベル1放送メモリ) に接続。
- L1BM から読み出されたデータは放送される (各 PE/MAB での書き込み制御で単一の PE/MAB に送ることもできる)。
- 分配モードも一応ある (あんまり速くない)

- 各 PE から読み出したデータは、縮約して L1BM に書き込むことができる。単一の PE/MAB を指定した読出しもできる。分配の逆操作 (結合) も一応ある。
- PE 間の直接結合はない。

## L1Bの特色

- 明示的な放送/縮約モードがあるため、複数 **PE** を使った細粒度並列処理を低オーバーヘッドで実行できる。
- **GPU** が非常に苦手とする **reduction** をオーバーヘッドを気にしないで使うことができる。
- 例えば比較的小さな行列乗算を **MAB** に分散させることができる。総和が遅いと、行列行列積  $A*B$  で **B** を縦に切る並列化をすると、各 **MAB** で **A** は全体が必要になる。総和が速いと、**B** を横に切る並列化ができ、**A** を分散させられる。また、単一の行列ベクトル積の並列化ができる。推論 (**LLM** でも) で非常に効率をあげやすくなる。
- (**PE** 命令と同期した、4 サイクル毎のデータ転送命令で制御)

## L2B(レベル2放送ブロック)

- **L1B** 8 個に対して **L2BM**(レベル2放送メモリ) が1つ。
- **L1B** 内と同様な放送・縮約・分配・結合・個別読出しを **L1B** に対して行う。
- データ幅は広くしている。
- いくつかの **L1B** 同士の直接データ交換モードがある。
- **PE** 命令と同期した、4 サイクル毎のデータ転送命令で制御。

# 全体

チップ4個を全体としてみると

- **L2B 16個、DRAM インターフェース4個、PCIe インターフェース4個**がある。
- **PCIe インターフェースはバッファメモリ (PDM) を持つ**
- **PDM-DRAM、PDM-L2B、DRAM-L2B** で、放送・縮約を伴うデータ転送を実行可能。
- このレベルのデータ転送命令は非同期で実行される。**PE 命令**とはお互いにタグ待ちができる。

## PE 命令の概要

PE 命令は

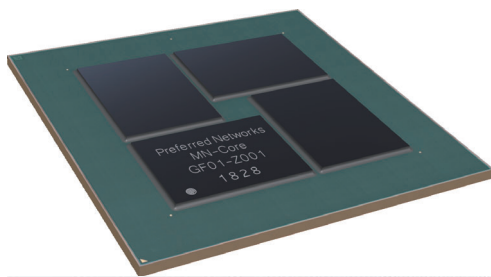
- 各演算の動作モード、入力セレクト、各メモリユニットのアドレス、**LM** では **read/write** モード、入力セレクト、各メモリユニットの入力マスクレジスタ番号、書き込みマスクレジスタ番号等
- **L1BM-PE, L2BM-L1BM** の転送モード、転送アドレス等

を指定する。

# 通常のアーキテクチャとの違い

- 現在の通常の **CPU** では、レジスタリネーミングと **OoO** 実行によって、最内側ループに本来ある並列性を実行時に回復する必要がある。これは、通常はアーキテクチャレジスタの数が不足し、コンパイラによるアンロールだけではパイプラインを埋められないからである(「京」は例外)
- **MN-Core** ではアーキテクチャレジスタ数の制限がないため、レジスタリネーミングが不要となる。
- さらに、固定長ベクトル命令の採用により、**OoO** 実行自体が不要になっている。直前の命令の実行結果が必ず利用可能だからである。(SVE や RISC-V vector extension でも実装可能)
- かなり単位の大きな **SIMD** 実行をすることで、非常に長い命令語を許しており、レジスタのアドレス長等の制限が不要になっている。

## MN-Core/MN-3 system



### PFNのスパコン「MN-3」が世界1位に、消費電力性能ランキングのGreen500で

2020.06.23

PFNのスパコン「MN-3」が世界1位に、消費電力性能ランキングのGreen500で

Preferred Networks (PFN) のスーパーコンピュータ「MN-3」が2020年6月22日(欧州時間)、スーパーコンピュータの消費電力性能ランキング「Green500」で世界1位を獲得した。HPC(ハイパフォーマンス・コンピューティング)に関する国際会議「ISC 2020 Digital」が同日ランキングを発表した。



PFNのスーパーコンピュータ「MN-3」

MN-3はPFNが独自開発した深層学習専用プロセッサ「MN-Core」を使ったスーパーコンピュータだ。PFNのスーパーコンピュータ「MN-2」の後継機で、2020年5月に運用を始めた。160個のMN-Coreを搭載し、1ワット当たり21.11ギ

## MN-Core 2 と後継

- **MN-Core2** 2023 年完成。販売開始した。
- **TSMC 7nm**。比較的小さいチップで **400TF** と **MN-Core** 並みの性能を実現
- 後継の開発もスタート
  - **Samsung 2nm**。完成時点で世界最高レベルの性能を目標
  - **LLM** 推論向けの開発も開始
- ポスト富岳向け検討もしてきたが、直近のものには採用されていない。

## MN-Core 2 のソフトウェア

- **MNSDK**: 機械学習向け
  - **PyTorch** 記述から **ONNX** 経由でコード生成
  - 既存の **PyTorch** コードから小さい修正で動作。ネットワークの記述等は同じ
- **HPCSDK**: 汎用 **HPC** 向け
  - **OpneCL** 方言(メモリモデルが違う)、**OpenACC** 方言を用意

## MN-Core 2 のアプリケーション性能

	MN-Core 2	A100
GCN(PFN internal use)	5.41TF(FP32)	3.17TF
ResNet50 training	77TF(FP16)	33.2TF(BF16)
ResNet50 Inference	154TF(FP16)	33.7TF(BF16)
HIMENO benchmark	9.03TF(FP32)	0.634TF
OpenFDTD	0.655TF(FP32)	0.488TF

- AI ワークロードでは同程度のピーク性能の **A100** の **1.5-5** 倍の性能
- 差分法系のアプリケーションやベンチマークでも高い性能。但し、**OpenFDTD** は **temporal blocking** を実装して **DRAM** アクセスを減らした成果。姫野ベンチはオンチップメモリにはいるサイズ。

## MN-Core の特徴

- チップ内の演算コアがそれぞれメモリをもち、コア間がツリーネットワークでつながる「チップ内分散メモリアーキテクチャ」
- キャッシュはない。外部 **DRAM** は共有される形で接続されるが、各コアがアドレスをだすのではなくデータ移動命令でツリーノードとの間でブロック転送
- チップ内分散メモリアーキテクチャと従来のアーキテクチャの中間
- 従来型プロセッサより高い性能と実アプリケーションでの高い実行効率を実現
  - **SIMD** 動作と、総和、放送をサポートする低レイテンシチップ内ネットワークによるオーバーヘッドの小さい並列化
  - 大きな **PE** ローカルメモリによる効率のよい **DRAM** アクセスの低減
  - 高い **PE** ローカルメモバンド幅 (ベクトル演算に対して **B/F=8**) によるカーネル部分の高効率実行

# 真のオンチップ分散メモリは？

- HBM は DRAM の平面的実装としてはほぼ究極。プロセッサダイのすぐ隣に DRAM ダイがある。
- それでも、電力消費はプロセッサダイと DRAM ダイの中を横に動く分、下のインターポーザの中を横に動く分が大きい。また、どうしても共有メモリアーキテクチャになる。
- 分散メモリにするは、DRAM ダイをプロセッサダイの上ののせて、さらに面的な配線で演算器とその直上(下?)の DRAM メモリブロックを接続すればよい。

## 現状の例: AP Memory/Powerchip

**WoW Hybrid Bonding**
                         
  **CoW uBump Bonding**



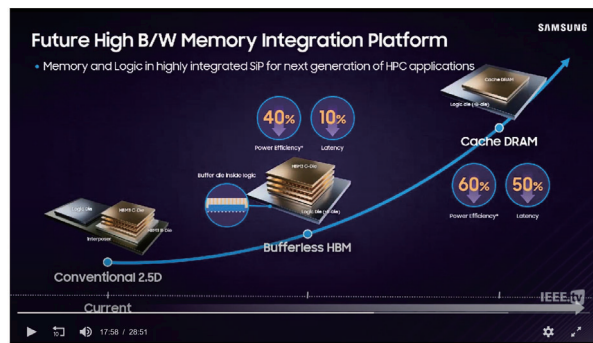
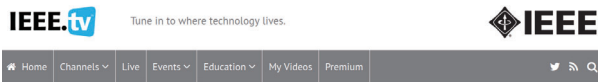
Stacking partner	Front End	Back End
pitch	Pitch < 3um >>100K IO	Pitch ~ 40um 1-2K IO
Base technology	BSI CIS	HBM
Application	Supercomputing Mining Networking	CPU/GPU AI AR/VR

**WoW Hybrid Bonding** 自体は既に実績がある技術。この会社は現状 **mining chip** を作っている模様

0.7pJ/bit 程度

<https://www.apmemory.com/> AP Memory 社のウェブサイトからハイブリッドボンディングのピッチ:現状でも 10 $\mu$ m。近い将来で 3 $\mu$ m とか=1 平方ミリで 10 万本。

# 将来の例



Heterogeneous Integration Platform for Next Generation Computing ('Beyond Moore')  
Published on February 24, 2023

Samsung の例: HBM に対して  
40-60% の電力削減

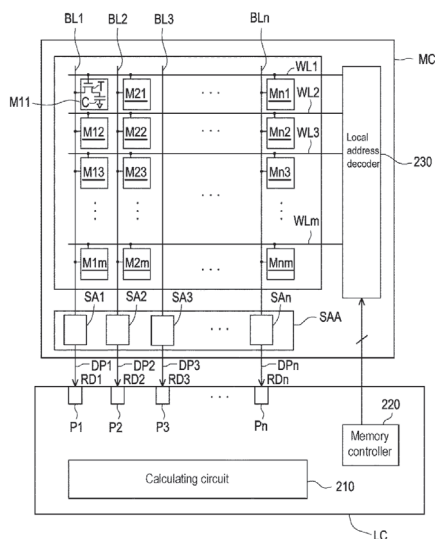
Cache DRAM で 1pJ/bit 程度  
(現状の HBMx は 3-4pJ/bit 程度)

Nanya、Winbond は PowerChip  
と同様な提案も (DRAM 6 社のうち  
下位 3 社)

## 原理的な限界は？

- HBM 等でも、DRAM 側の現状の設計は DDRx の DRAM と同等。クロックがはいり、FIFO やセレクタを通してデータがでてくる。部分書き込みのためのデータ保持回路等もある。
- DRAM チップ内でみても、実際の DRAM セルやセンスアンプの消費電力は結構小さい (O'Connor et al. Fine-Grained DRAM: Energy-Efficient DRAM for Extreme Bandwidth Systems, MICRO-50)
- ハイブリッドボンディングでは非常に大きなパッド数が可能であるので、回路構成から再検討が可能。
- 0.1-0.2 pJ/bit くらいまではそのうちいけそう。

# PSMC 特許(出願中)



## PSMC 特許 US20230186972A1 Memory device with high data bandwidth

これは多分一番極端な回路構成案。センスアンプから先をすべてプロセッサダイ側にもってくる。

(特許申請自体は **DRAM** をロジックダイにのせますと書いてあるだけなのでこのまま成立はしなさそう。それは前例がある)

## 実装上の懸案事項

一杯ある。

- **HBM** みたいに **DRAM** スタックできるし、**TSV middle** とかですごく沢山パッドだせるんだけど、**DRAM** スタック上につんで冷却できるか？
- **face up** でロジックダイ使って電源とか信号線引き出せるか？
- **DRAM** スタック下においたら、**DRAM** スタック通して信号とか電源だせるか？

とはいえ、解決できない問題という感じでもない。物理的に金属配線でつながるので微細化できれば数増えてバンド幅あがる。(磁界結論とか容量結合はスケールしない)

# MN-Core アーキテクチャと 3D DRAM

- **MN-Core** アーキテクチャはチップ全体が **SIMD** アレイ。それぞれの **PE** がローカルメモリをもって、階層的なネットワークで結合。
- **3D DRAM** は、基本的には各 **PE** に **DRAM** ブロックがつく。例えば **16MB** とか。
- **80** 年代風の超並列 **SIMD** プロセッサとほぼ同じ。ネットワークはさぼっている。概ね同じようにプログラミングもできる。
- つまり: ほぼ **HPF (OpenACC** ともいう) でプログラムできる。(但し、**DRAM** とオンチップメモリの違いを意識する必要はある)

## チップ内分散メモリでアプリケーションは書けるか？

- 概念的には、**1990** 年代に共有メモリベクトルプロセッサから分散メモリマシンに移行したのとかわからない。
- 数値風洞 (**VPP500**)、**IBM SP-2**、**Cray T3D**、**PC** クラスタ。 **MPI** で書いてね的、、、
- 但し、**MN-Core** の場合 **SIMD** なので、そのもうちょっと前の **SIMD** 超並列 (**TMC CM**, **Maspar MP**) にアーキテクチャ、プログラミングモデルは近い。**CM-Fortran (HPF** の原型) みたいなものが見える。比較的ネットワークが高速。

# 他社の AI むけプロセッサは？

- 今年の HotChips 24: 口頭発表 24 個のうち 14 個が AI アクセラレータ。
- Tenstorrent Blackhole、SK Hynix PiM、Blackwell、Sambanova SN40L、Intel Gaudi 3、AMD MI300X、FuriosaAI RNGD、AMD(Xilinx) Versal、Onyx (Stanford の研究)、Meta MTIA、Tesla Dojo、Cerebras CS-2、MS MAIA、PFN MN-Core 2
- Blackhole、SN40L、Onyx、MTIA、Dojo、Cerebras、MAIA: 2D オンチップネットワークに MIMD コア+行列演算ユニット。
- Gaudi、RNGD: 大きな行列乗算器もつ独自アーキテクチャ(Groq も)
- SK Hynix は GDDRx メモリチップの中に演算器いれるもの
- Versal は FPGA、あとは Nvidia, AMD, PFN

## アーキテクチャの特徴

- 階層キャッシュはもうないというのは Nvidia と AMD の GPU 以外では共通。
- AI プロセッサアーキテクチャはほとんどが MIMD 2D メッシュ。メッシュの辺のところに HBM つける。
- どうやってプログラムを書くのか?というとりあえずとても大変そう。
- AI 用プロセッサは FP64 ない。FP32 も絶滅しつつある。Nvidia、AMD の GPU も FP64 の性能が相対的に下がってきている。(MN-Core は異なる精度間で演算器を共有しているので比較的 FP32, FP64 の性能を維持しやすい)

## まとめ

- 計算機の進歩について、「ポストムーア」は今度は本当。シリコン半導体技術の限界は確実にきている。
- しかし、計算機設計という観点から見ると、問題はそこではなくて、メニーコア **CPU** や **GPU** の設計が半導体技術の進歩を上手く利用できなくなっていること。これは **30** 年前にシステムレベルでの共有メモリ計算機が無理になってきたのと同じ。
- **GPU** アーキテクチャは階層キャッシュ構造から離れるのは難しいように見える。
- **AI** 向けプロセッサはキャッシュなし、ローカルメモリの方向。
- **MN-Core** は、さらに大規模 **SIMD** と階層的オンチップネットワークで高いピーク性能と **AI** でも汎用 **HPC** でも高い実行効率を実現。
- さらに、**3D DRAM** を有効に利用して **DRAM** バンド幅をあげることが重要になる。我々はこれにも取り組んできている。



アドバンスソフト  
からの情報提供

1. 今後のセミナー予定
2. セミナー資料のダウンロード
3. アーカイブ動画のご案内

Copyright ©2024 AdvanceSoft Corporation. All rights reserved.

## 今後のセミナー予定

第7回

2024年10月11日（金）  
14：00～15：30

「建築・都市の新たな風環境デザイン」

千葉工業大学 創造工学部 建築学科  
元教授 森川 泰成 様

受付中

第8回

2024年11月22日（木）  
14：00～15：30

「極低温超伝導コンピュータ  
アーキテクチャ技術とその新展開」

九州大学 大学院システム情報科学研究院  
情報知能工学部門 主幹教授 井上 弘士 様

受付中

Copyright ©2024 AdvanceSoft Corporation. All rights reserved.

下記のURLから、過去のセミナー資料をダウンロード可能です。  
<https://www.advancesoft.jp/download/>



ホーム / シミュレーション図書館 (資料ダウンロード)

### シミュレーション図書館 (資料ダウンロード) について

- 初めてご利用いただくお客さまは、「新規メンバー登録 (フォーラム会員登録)」をお願いいたします。ご登録いただきますと、いただいたメールアドレスあてにダウンロード用のログインパスワードをお知らせいたします。
- ログインパスワードをお持ちのお客さまは、お客さまのご登録されたメールアドレスと、こちらからお知らせしましたパスワードで、ログインしてください。ログイン中ほどの資料もダウンロードいただけます。
- 会社案内、会社概要、開発実績は、ログイン不要でダウンロードいただけます。
- ご登録いただいたメールアドレスあてに、メールマガジンやアドバンスソフトからのご案内をお送りさせていただきます場合がございますので、ご了承ください。
- 資料に含まれている内容を、その一部でも著作権者の許諾なしに、複製、改変、配布を行うことおよびインターネット上で提供する等により、一般へ送ることは法律によって固く禁止されています。

### 資料カテゴリー一覧

パッケージソフトウェア▷ ナノ・バイオ 半導体 流体 管路系 二相流 防災 構造 統合環境  
解析・コンサルティング▷ 防災・地震 二相流 2次電池 その他



ホーム / シミュレーション図書館 (資料ダウンロード) / パッケージソフトウェア / 半導体

### 半導体

【半導体】> セミナー資料 (アドバンスソフトの半導体デバイス、電磁波解析のご紹介 ~DX、グリーンイノベーションを加速する数値解析シミュレーション~) 20221117\_all 【全体】

カテゴリー [半導体](#) [パッケージソフトウェア](#)

【半導体】> セミナー資料 (半導体デバイス3次元TCADシステム Advance/TCAD 最新動向セミナー) 20181225\_all 【全体】

カテゴリー [半導体](#) [パッケージソフトウェア](#)

【半導体】> セミナー資料 (半導体デバイス3次元TCADシステム Advance/TCAD 製品紹介セミナー) 20190828\_all 【全体】

カテゴリー [半導体](#) [パッケージソフトウェア](#)

【半導体】 Advance/TCAD > セミナー資料 (オンラインセミナー「アドバンスソフトの電磁波、半導体デバイス解析のご紹介 ~IoTを支える数値解析シミュレーション~」) 20200825\_all 【全体】

カテゴリー [半導体](#) [パッケージソフトウェア](#)

【半導体】 Advance/TCAD > セミナー資料 (半導体デバイス3次元TCADシステム発売のご紹介セミナー)

Copyright ©2024 AdvanceSoft Corporation. All rights reserved.

YouTubeにて、セミナーのアーカイブ動画を配信中  
<https://www.youtube.com/user/advancesoft>



- 【セミナー動画】第1回 アドバンス・シミュレーション・セミナー 2024年4月19日 (金) 開催『ロケットと宇宙機に関するいろいろな自由研究』 advancesoft • 235 回視聴 • 12日前
- 【セミナー動画】第8回 アドバンス・シミュレーション・セミナー 2023年10月6日 (金) 開催『爆轟から見える CAE の方向性』青山... advancesoft • 337 回視聴 • 6か月前
- 【セミナー動画】第6回 アドバンス・シミュレーション・セミナー 2023 2023年8月28日 (月) 開催『フェイクニュースと生成 AI』... advancesoft • 601 回視聴 • 8か月前
- 【セミナー動画】第5回 アドバンス・シミュレーション・セミナー 2023 2023年8月3日 (木) 開催『原子力安全に必要な計算... advancesoft • 343 回視聴 • 9か月前



- 【セミナー動画】 NanoLabo/NeuralMD 新機能および今後の開発計画のご紹介 2023年10月18日 (水) 開催 advancesoft • 231 回視聴 • 6か月前
- 【セミナー動画】 生成 AI・大規模言語モデルとCAE セミナー 2023年9月22日 (金) 開催 advancesoft • 953 回視聴 • 7か月前
- 【セミナー動画】 NanoLabo/NeuralMD 新機能紹介セミナー 2023年4月20日 (木) 開催 advancesoft • 385 回視聴 • 1年前
- 【セミナー動画】 OpenModelica 活用セミナー 2022年11月30日 (水) 開催 advancesoft • 838 回視聴 • 1年前

Copyright ©2024 AdvanceSoft Corporation. All rights reserved.





**警告**

このレポートに収録されている文章および内容については、ご自身のために役立つ用途に限定して無料配布しています。  
このレポートを、販売、オークション、その他の目的で利用するには、著作権者の許諾が必要になります。  
このレポートに含まれている内容を、その一部でも著作権者の許諾なしに、複製、改変、配布を行うことおよびインターネット上で提供する等により、一般へ送ることは法律によって固く禁止されています。