

AIコンピューティング研究ユニット

本村研 藤木研

紹介資料

26/3/3

総合研究院

AIコンピューティング研究ユニット (ArtIC)

情報通信系 藤木研

ArtICの成り立ち

情報通信系

本村・劉研究室 => 本村・藤木研究室
(23年2月まで) (23年11月から)



=
(イコール)

総合研究院

AIコンピューティング研究ユニット

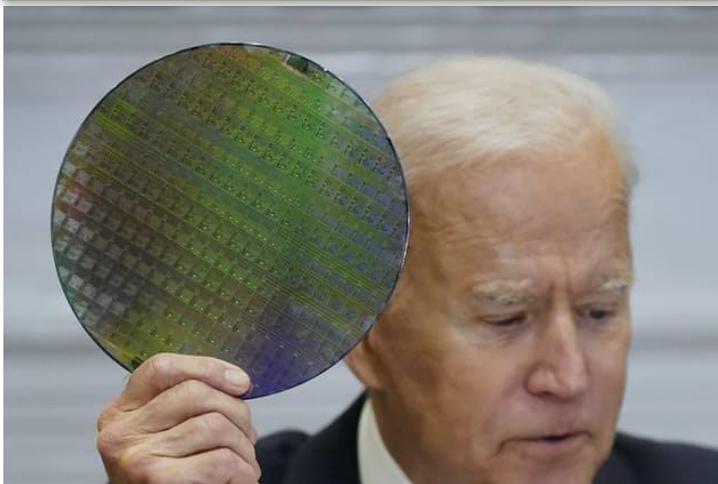


Artificially Intelligent Computing Research Unit

もう一つの意味: **ArtなIC(集積回路)** => ソフトとハードの総合技術・芸術

ArtIC 5年間の社会情勢の変化

バイデン大統領 半導体 \$50B 投資宣言 (2021)



TSMC熊本工場(2022-)



ラピダス千歳工場(2023-)



Chip War (2022)



AI時代においては、データこそが新たな石油だとよく言われる。しかし、私たちが直面している真の制約は、データではなく処理能力の不足にある。データの保存や処理ができる半導体の数は限られており、その製造工程は目が回るほど複雑で、恐ろしいまでのコストがかかる。いろいろな国から購入できる石油とはちがって、計算能力の生産はいくつかの決定的な急所にまるまるかかっている。それは、一握りの企業、ときにはたった1社でしか生産できない装置、化学薬品、ソフトウェアである。(序章より)

定価2970円(本体2700円+税10%)

AIの勃興と国際情勢の緊迫化



半導体・集積回路・チップ(すなわちハードウェア)が、技術的にも社会的にも注目の的へ

AIコンピューティング研究の背景

急ピッチで進むAI技術の社会応用



自然言語・自動
対話型AI



スマート
ロボット



自律航行
ドローン



スマート社会
インフラ

AI処理の急拡大 => エネルギー消費問題の深刻化

より豊かで低環境負荷なスマート社会を実現したい



「AIコンピューティング」の処理効率を向上する
情報処理アーキテクチャの技術革新が必要

ChatGPT Burns Millions Every Day. Can Computer Scientists Make AI One Million Times More Efficient?

John Koetsier Senior Contributor @
John Koetsier is a journalist, analyst, author, and speaker.

Follow

Feb 10, 2023, 03:09pm EST

Listen to article 8 minutes



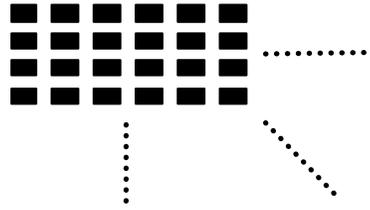
Running ChatGPT costs millions of dollars a day, which is why OpenAI, the company behind the viral natural-language processing artificial intelligence has started ChatGPT Plus, a \$20/month subscription plan. But our brains are a million times more efficient than the GPUs, CPUs, and memory that make up ChatGPT's cloud hardware. And neuromorphic computing researchers are working hard to make the miracles that big server farms in the clouds can do today much simpler and cheaper, bringing them down to the small devices in our hands, our homes, our hospitals, and our workplaces.

(Forbes)

ChatGPTの学習には
一般家庭の数百年分も
の消費電力が必要
(Stanford大学報告書)

一つの視点: 大脳とGPU

GPU(数万モジュール)



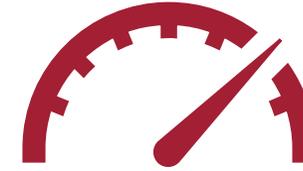
演算性能

数 Exa
(10^6 Tera)
Ops/sec

メモリ容量

Peta
(10^6 Giga)
Byte

消費電力



数MW

GPU(単体モジュール)



演算性能

100 Tera
Ops/sec

メモリ容量
(iPhone並み)

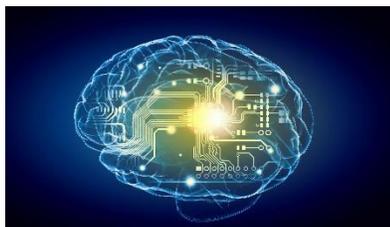
32
Giga
Byte

消費電力



400W超

人間の脳



2 Tera
Ops
/sec

演算性能
(iPhone並み)

Peta
(10^6 Giga)
Byte

メモリ容量

~20W



消費電力

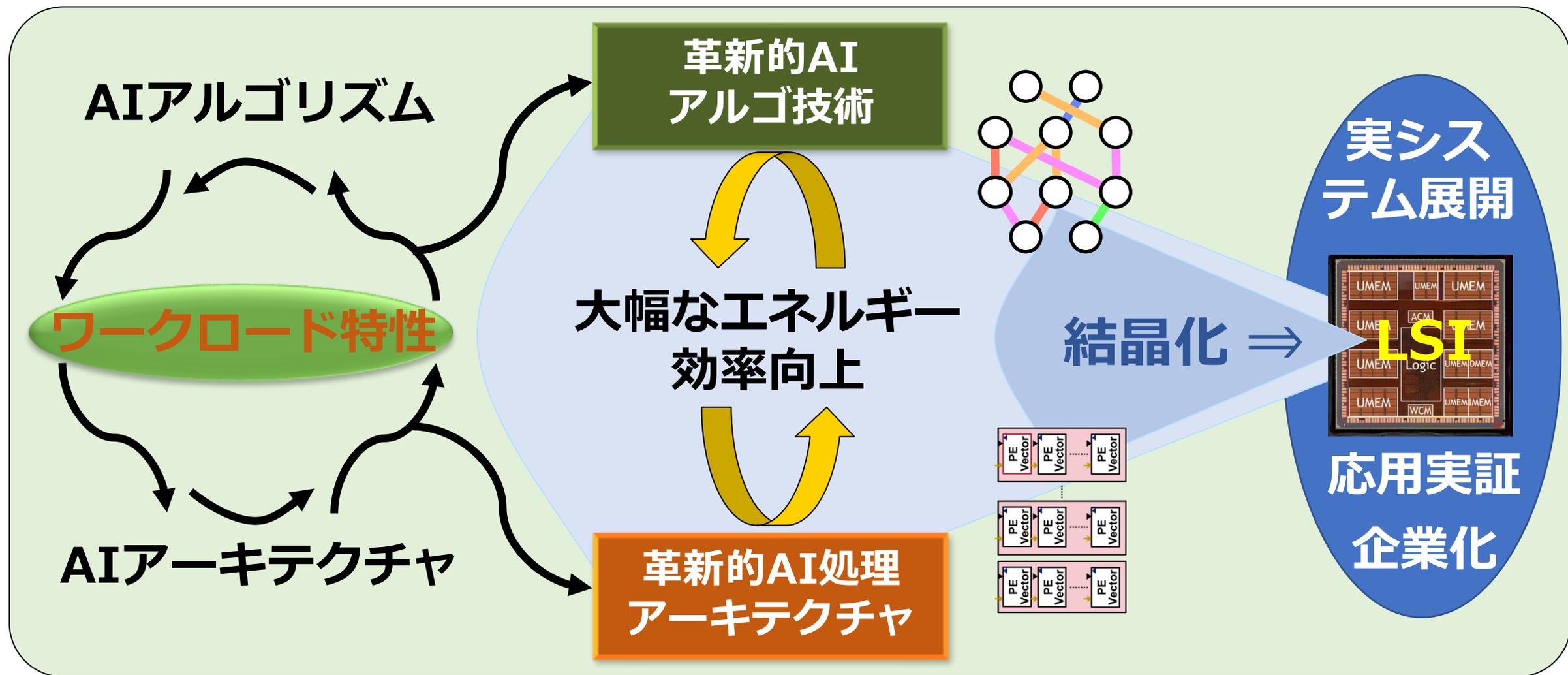


工夫すべきところは
山ほどある



研究しよう!

ArtICの研究スタイル



SW-HWの両面にわたって研究を推進

ArtICの研究スタイル



AIアルゴリズム



システム・SW



EDGE
COMPUTING



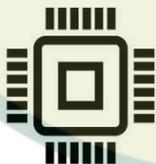
データセンタ



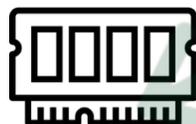
計算機アーキテクチャ



GPU



アクセラレータ



メモリセントリック計算

エッジ計算からサーバーまで様々なフォームファクタを俯瞰し、1チップの開発に閉じない、効率と柔軟性を両立する次世代の計算機基盤を実現

- ・データ移動コストに注目したメモリ内計算技術とそのシステム応用
- ・AI向けメモリ最適化コプロセッサ
- ・HW/SWコデザインを通した高効率システム

大規模生成モデル・ゲノムプロセッシング・プライバシープリザード計算・リアルタイムデータベースなどの重要アプリの効率化へ

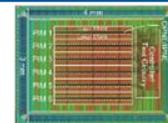
SW-HWの両面にわたって研究を推進

ArtICの世界最先端AIチップ実績

半導体
プロセス

- 二値化DNNアクセラレータチップ
■ **VLSIシンポジウム2017**
- 対数量子化・三次元積層DNNアクセラレータチップ
■ **ISSCC2018, JSSC2018**
- 完全結合型・全並列型デジタルアニーリングチップ
■ **ISSCC2020, JSSC2020**
- シフト演算型・直積アレイ型DNNアクセラレータチップ
■ **Hot Chips 2021**
- ランダム重み固定型DNNアクセラレータチップ
■ **ISSCC2022, ICML2022**
- 完全結合型メタモルフィックアニーリングチップ
■ **ISSCC2023**
- 超低電力・ビット累進型DNNアクセラレータチップ
■ **VLSIシンポジウム2023**
- 三重・不規則スパース性対応DNNアクセラレータチップ
■ **A-SSCC2024**

太字下線部は
世界初の技術



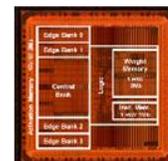
65nm



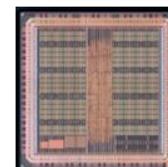
40nm



65nm



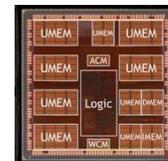
40nm



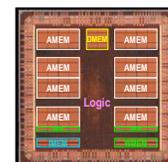
40nm



40nm



40nm



40nm

最先端のAIアルゴリズムに注目しアーキテクチャ・イノベーションにこだわって研究を推進。

目玉となる成果に絞ってチップ実装・実機デモ。

チップ化は成果の極く一部であり、アルゴリズム-アーキテクチャ融合型の研究が、研究活動の大半。

ArtICの世界最先端アーキテクチャ研究実績

- In-ReRAMベクタ計算アーキテクチャ
 - ASPLOS 2018
- オートマトン型ゲノムプロセッサアーキテクチャ
 - ISCA 2018
- In-SRAM SIMT計算アーキテクチャ
 - ISCA 2019
- DP型ゲノムプロセッサアーキテクチャ
 - MICRO 2020
- 複数階層型インメモリ計算アーキテクチャ
 - MICRO 2022
- ビュー生成によるコヒーレントPIMアーキテクチャ
 - MICRO 2023
- モバイル向けIn-SRAM PIM ISA
 - HPCA 2025
- GNNアクセラレータアーキテクチャ
 - ISCA 2025
- PIM向けLLM/Attention圧縮手法
 - HPCA 2026

データ移動コストを大幅に低減するPIM/IMC技術に注目しアーキテクチャ・イノベーションにこだわって研究を推進。

コンパイラやシステム設計など、計算機スタックを縦断的に分析し実装・及びシミュレーション評価。

ゲノム分析などのAI以外の重要アプリにも着目し、アルゴリズム-アーキテクチャ融合型の研究を実施。

ArtIC人員構成

		本村研	藤木研	
教員	教授	本村 真人		全体運営、ディープラーニング
	准教授		藤木 大地	アーキテクチャ
	助教	金子 竜也		
	特任助教/ポスドク	大越 (9月まで/Imperialへ訪問研究員として出向)		
秘書		橋本, 土屋		
学生	博士	6名 (9名→3名卒業、昨年M2は0名)		
	修士	8名 + 募集若干名		
	学部	2名 (昨年度から藤木研のみ)		



- 機械学習アルゴリズムや, 計算機アーキテクチャ, ハードウェア設計などに興味を持つ皆さんのArtICへの参加を歓迎します
- 二研究室で共同運営しており, 研究ユニット内に垣根はありません
 - どちらを志望しても全く差はありません
 - ArtIC教員全員で協力して, 丁寧な指導と居心地よい環境づくりとを心掛けています

研究対象の技術レイヤ

技術レイヤ

「論理」を相手にする世界



「物理」を相手にする世界

身に付くスキル・知識 (例)

- 機械学習 アルゴリズム
- 画像処理 アルゴリズム
- ディープラーニング(人工知能)
- 学習環境構築・GPU活用
- 組み込みシステム設計
- プロセッサアーキテクチャ
- 非線形・近似計算理論・応用
- 柔軟いHW (動的再構成LSI)
- 信号処理
- デジタル回路・Verilog設計
- FPGA設計・利用
- LSI設計

SW-HWの両方に興味がある人向き

- 研究に集中できるように、学部4年生から**全員RA**雇用しています
- 研究環境(オフィス, 計算機, 実験評価)の整備には力を入れています
- 導入教育 (~3月)
 - 輪講: アーキテクチャ, 機械学習
 - 実習: ディープラーニング, FPGA
- 研究テーマ配属
 - 4年生 4月頃。本人の希望に応じてテーマ調整
- 研究の進め方
 - 全体会議: 週1回 (バーチャル)
 - 4グループ毎の研究報告: 週1回 (バーチャル)
 - その他, 適宜個別に打合せ・議論 (対面)
 - コアタイムは設けていません
- **修士の6割(過去実績)が博士進学**、4割が企業就職
博士の大半は企業研究職に就職

まだ始まったばかりの研究ユニット(研究室)です。自由に闊達な雰囲気、創造的なアイデアを生み出せる環境づくりを目指しています。

真新しい建物、美しく広く眺望の良いオフィス(すずかけ台・J3棟17階フロア全体)という恵まれた環境を生かし、魅力的な居室環境の整備を進めています。構成メンバーの意見を柔軟に取り入れながら、居心地が良く研究モチベーションが湧いて出るような生活環境を目指します。

所属学生には、広めのデスクスペース、ノートPCと32型ディスプレイを支給します。深層学習用計算サーバやFPGAボード等、研究環境も充実しています。共同研究資金・競争的研究資金による博士・修士学生のRA雇用も積極的に進めています。



修論・卒論 タイトル名 (2025年度)

□ 修士論文

- N/A

□ 卒業論文

- パラメーター凍結LLMのロスレス圧縮手法の検討、岩田
- CPUによるデータの局所性を高めた近似計算とPIMによるその精度補償、多田
- アクションのキャッシングと再拡散によるVision-Language-Actionモデル推論の高速化、大井

修論・卒論 タイトル名 (2024年度)

□ 修士論文

- 制約付き組合せ最適化問題の解探索を高効率に行うアニーリングプロセッサ、兵藤
- 説明の後方互換性を考慮した勾配ブースティング決定木の再訓練法、山倉
- 強い宝くじ仮説とグラフ分割を用いたGNN 処理効率化の研究、伊藤
- 乱数部分ネットワーク内におけるコンパクトな強い宝くじの研究、大塚
- 平均場近似の高精度化に基づく高性能な二次無制約二値最適化手法、黒木
- エッジデバイス向けのニューラルネットワーク圧縮手法、塩田

□ 卒業論文

- 鍵共有不要な可変長データPIR、中森
- PIMによる動的量子化を用いた大規模言語モデル推論の効率化、松島
- ニューラルネットワークのエッジ向け継続学習手法、石橋

修論タイトル名 (2023年度)

□ 修士論文

- A Highly Accurate and Parallel Vision MLP FPGA Accelerator based on FP7/8 SIMD Operations and efficient dataflow design、安永
- イジング計算機を対象とした動的パラメータ調整の研究、井上
- バケッティング・データ構造による自己位置推定機構のメモリ削減及び高速化、市川
- 負荷均等配分を目指した高並列疎行列積アーキテクチャの研究、永原
- 局所鋭敏性ハッシング機構を用いた超次元コンピューティングの研究、渡邊

修論・卒論 タイトル名 (2022年度)

□ 修士論文

- 全並列アニーリングの解探索性能を向上させる動的なスピン反転機構の研究、小此木
- 局所解脱出を容易にするアニーリング手法とそのアクセラレータ設計、神保
- 乱数重みニューラルネットワークにおける精度・サイズトレードオフの向上に関する研究、大越
- 高効率な量子化決定森推論アクセラレータのためのモデル最適化手法の研究、北島

□ 卒業論文

- 強い宝くじ仮説に基づく超軽量物体検出ネットワーク、大塚
- 同変性ネットワークに基づく自律走行向け強化学習手法、塩田
- 表形式データを対象とした決定木とニューラルネットワークの融合型機械学習手法の研究、山倉
- 2スピン同時フリップを並列試行可能なシミュレーテッド アニーリング手法、兵頭
- 組合せ最適化問題のアニーリング解法に関する難易度評価、四元

最後に…

- ArtICホームページを確認ください（東京科学大 ArtIC）
- 機械学習アルゴリズムや、計算機アーキテクチャ、ハードウェア設計に興味を持つ皆さんのArtICへの参加を歓迎します
 - 特別な知識は求めません。この分野の研究に対する意欲を期待します
 - 4年生前半までに基礎知識が身に付くよう、輪講や研修を行います
- 二研究室で共同運営しており、研究ユニット内に垣根はありません
 - どちらを志望しても全く差はありません
 - ArtIC教員全員で協力して、丁寧な指導と居心地よい環境づくりとを心掛けています
- 一線級の国際会議で発表できるグループです
- 実戦的な研究活動を主体としています
- 実社会で役立つ考え方・スキル・知識を身に着けることができます
- 産学連携、大学間連携、国家プロジェクト参画を活発に進めています





J2

J3







フロア案内

FLOOR INFORMATION

17F



1703	共通事務室1
情報工学系 Department of Computer Science	
山村研究室 Yamamura Lab.	
1706	教員室(山村) Professor M.Yamamura
1707	研究室
1710	学生室
小野 功研究室 Isao Ono Lab.	
1704	教員室(小野 功) Associate Professor Isao Ono
1705	学生室
1709	研究室
科学技術創成研究院 Institute of Innovative Research	
AIコンピューティング研究ユニット AI Computing Unit	
1711	ユニット控室
1712	学生室
1713	教員室(本村) Professor M.Motomura
1714	教員室(劉) Associate Professor J.Yu
1715	秘書室・会議室
1716	教員・研究員室



← 1711-1716

1713-1714-1715 →

715

ARTIC

ARTIC



押

← 1711-1716

1713-1714-1715 →

1715

ArtIC
秘書室

教授(本学)
Professor M. Hoshimura

准教授(本学)
Associate Professor J. Fu

秘書室 担当



