

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

25/11/12

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 本村研 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)

The screenshot shows the university's research website interface. It displays three research groups: Sasaki Research Group, Momota Research Group, and Nakamura Research Group. The Momota Research Group page is highlighted with a red border. The page title is "次世代AIコンピューティングアーキテクチャを創成する" (Creating the next-generation AI computing architecture). Below the title, there is a brief description of the research focus on AI computing architecture and hardware. The page also includes a photo of the researcher and links to his profile and publications.

(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

The screenshot shows the website for the AI Computing Unit. The header features the university's logo and the unit's name. The main content area includes a brief overview of the unit's mission, a portrait of the leader (Masato Momota), and a detailed profile section. The profile lists Momota's academic background, including his education at the University of Tokyo and his work at the University of California, Berkeley. The website also features a diagram illustrating the research focus on accelerating AI computing architecture.

2019年4月
に発足



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

本村・藤木の自己紹介

本村

テクノロジストの時代

世界の研究機関や大企業が力を籠った人材が、AIへの専用導体の研究開発のトップグループに、日本から食い込んでしまっているが東京工業大学の本村眞人教授が、その研究結果を次々と発表している。1988年に入社したNECの中央研究所で集積回路の研究に取り組んだのをきっかけに、半導体を研究する道を歩み始めた。専門分野のア

東京工業大教授
本村 真人氏

1987年京大理学部修了
NECにて社員登用試験で合格。同年同大博士号取得。
マサチューセッツ工科大学にて現職。59歳

AI半導体世界と競う

1987年京大工学博士
'96 NEC研究所
'11 北大
'19 東工大
ArtIC立ち上げ

LSIのオリンピック
ISSCCで次々に世界初のチップを発表

4千人参加の最高峰会議

多くの国内外シンポジウムでAIハードウェアの招待講演

日経新聞
21/10/26日朝刊

II 随時掲載

- '87 京大理学部修士
- '96 京大工学博士
- '87-'11 NEC研究所
- '11-'18 北大
- '19- 東工大

ArtIC立ち上げ

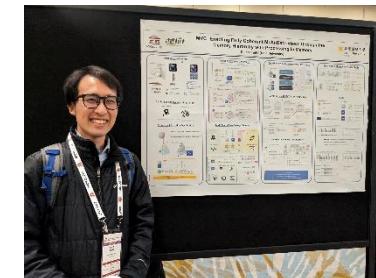


多くの国内外シンポジウムでAIハードウェアの招待講演

藤木

- '16 慶應理工学部卒
- '16 ミシガン大 Ph.D.課程入学
- '22 ミシガン大 Ph.D.取得
- '22-'23 慶應大学 助教
- '23.11月 東工大 准教授

ArtICに参加



計算機アーキテクチャの5大トップ国際会議で次々に論文発表。日本には稀有な存在

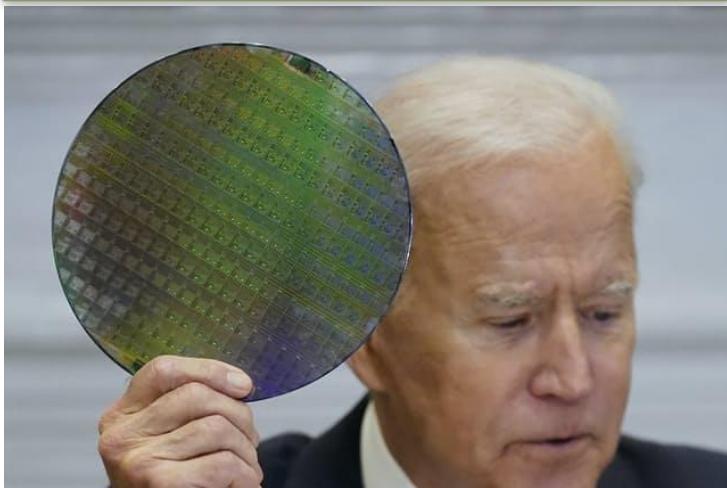


若くして既に著書も…

共同研究@Nvidia, Samsung, Meta FAIR

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(つなわち
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.

2

Listen to article 8 minutes



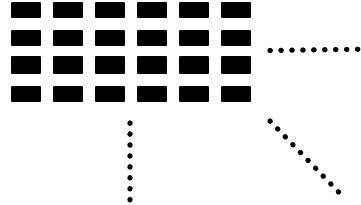
Running ChatGPT costs millions of dollars a day, which is why OpenAI, the company behind the viral natural-language processing artificial intelligence has started ChatGPT Plus, a \$20/month subscription plan. But our brains are a million times more efficient than the GPUs, CPUs, and memory that make 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)



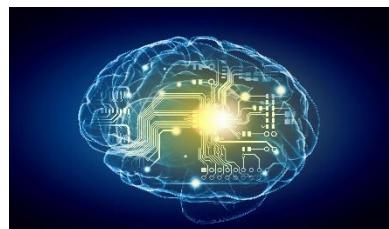
消費電力



GPU(単体モジュール)



人間の大脳



消費電力

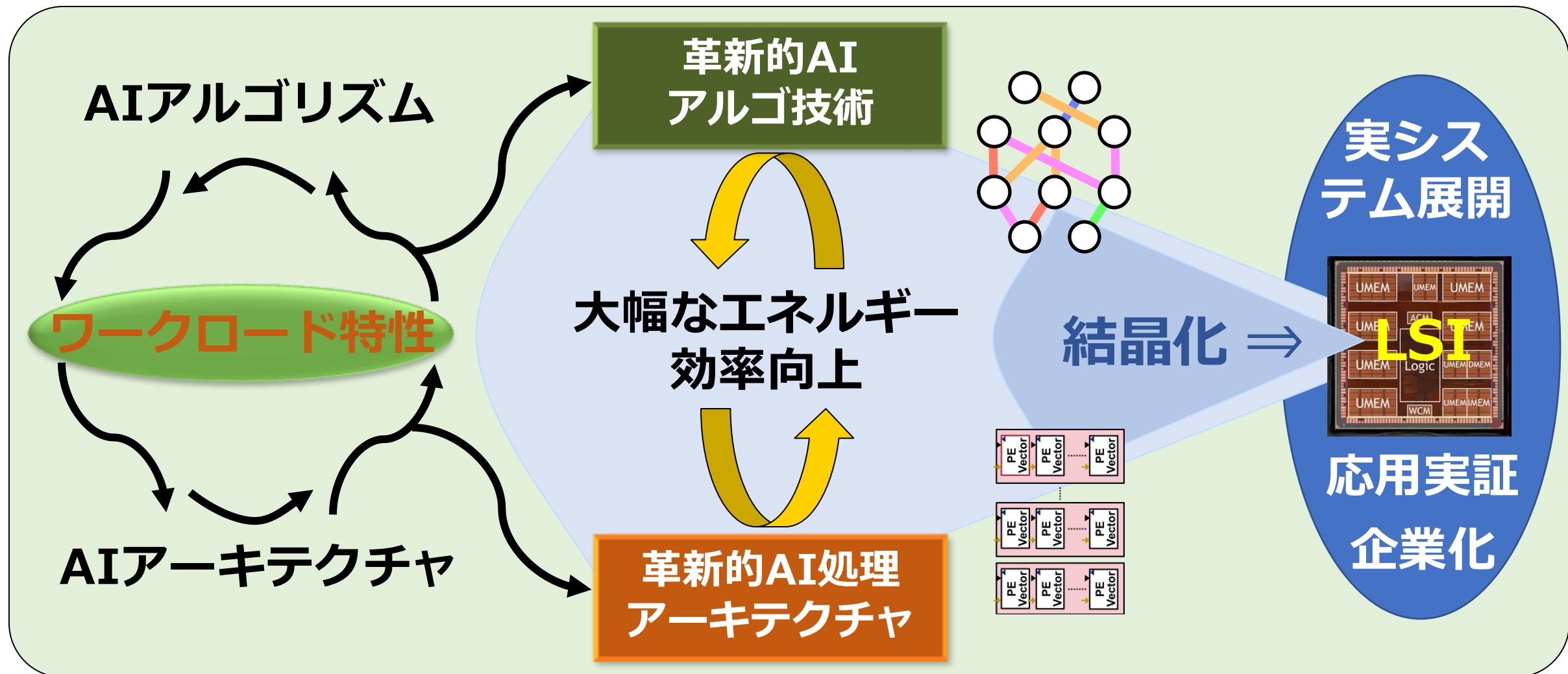


工夫すべきところは
山ほどある



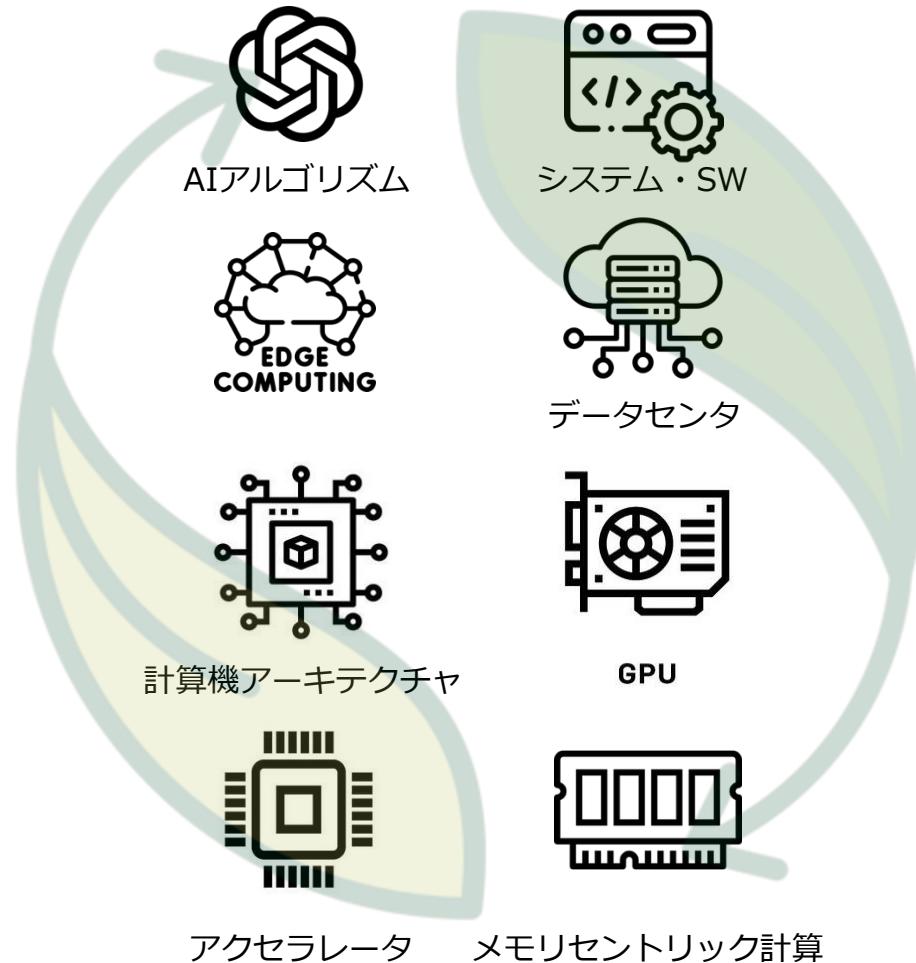
研究しよう！

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



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
- ・AIワークロードを意識した時空間展開型コンピューティング
- ・HWによる高効率実行に最適化したAIアルゴリズム

大規模生成モデル・ゲノムプロセッシング・プログラミング・プリザーブド計算・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

□ 二値化DNNアクセラレータチップ
■ **VLSIシンポジウム2017**

太字下線部は
世界初の技術

□ 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**

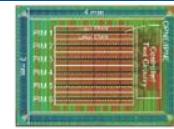
□ ランダム重み固定型DNNアクセラレータチップ
■ **ISSCC2022、ICML2022**

□ 完全結合型メタモルフィックアニーリングチップ
■ **ISSCC2023**

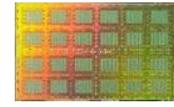
□ 超低電力・ビット累進型DNNアクセラレータチップ
■ **VLSIシンポジウム2023**

□ 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**

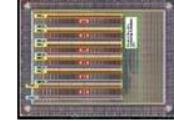
半導体
プロセス



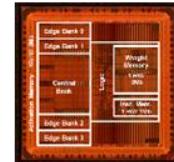
65nm



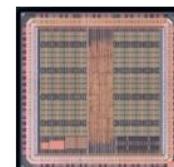
40nm



65nm



40nm



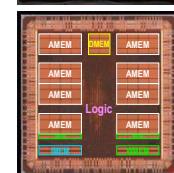
40nm



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 SIMD計算アーキテクチャ

■ **ISCA 2019**

■ GPU内SpMMデータ変換コプロセッサ

■ **SC 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

■ **ISCA 2025**

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以外の重要なアプリにも着目し、アルゴリズム-アーキテクチャ融合型の研究を実施。

Fixed-Random-Weight DNN Chip @ ISSCC2022

Algorithm

Background for Deep Learning

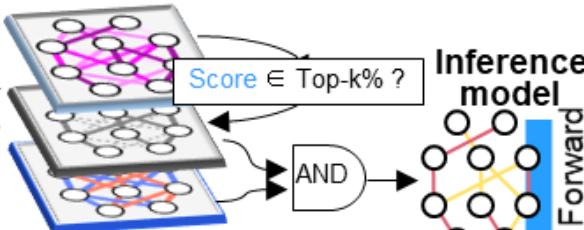
- Requirements for external memory access reduction
- The advent of a new efficient algorithm

Hidden Network (HNN) [V. Ramanujan+, CVPR2020]

Scores trained by backpropagation

Supermask for top- $k\%$ connections

Initial model with random **Weights**



Key Features

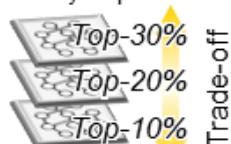


Architecture Implication

Weight: {-1, +1} **Supermask:**

$$\begin{cases} 1 & \text{--- } k\% \\ 0 & \text{--- } (100-k)\% \end{cases}$$

Exchanging only supermask



On-chip model construction

RNG

Seed

Generate Weights

No need to store

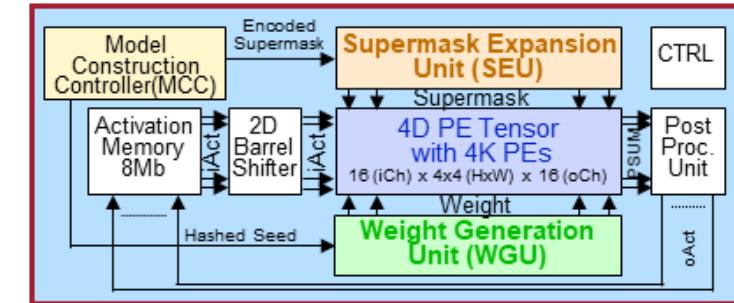
X

Expand Supermask

Zero-Run-Length (ZRL) compression

Architecture

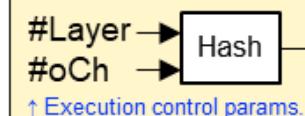
Hiddenite Chip



On-chip Weight Generation

Hashed seeds eliminate the need to store weights without accuracy degradation

MCC

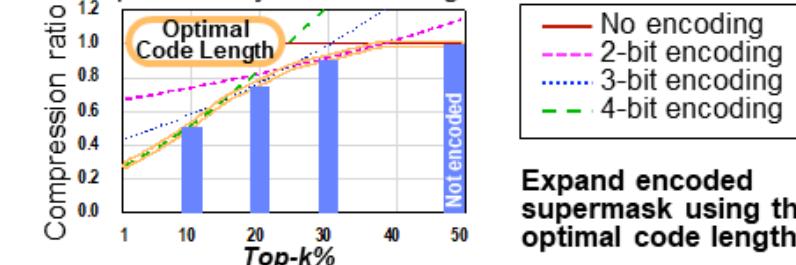


WGU



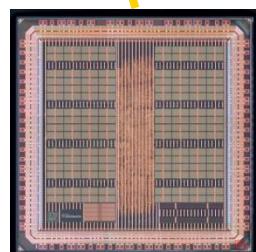
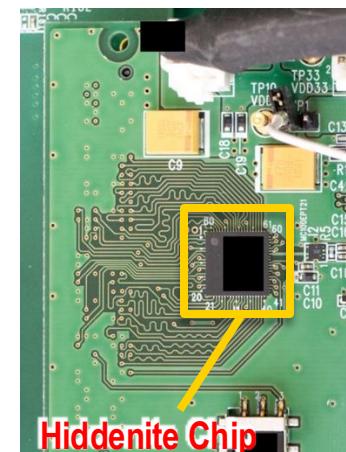
On-chip Supermask Expansion

Compression by ZRL encodings



Expand encoded supermask using the optimal code length

Real Chip & Demo System



ArtIC人員構成

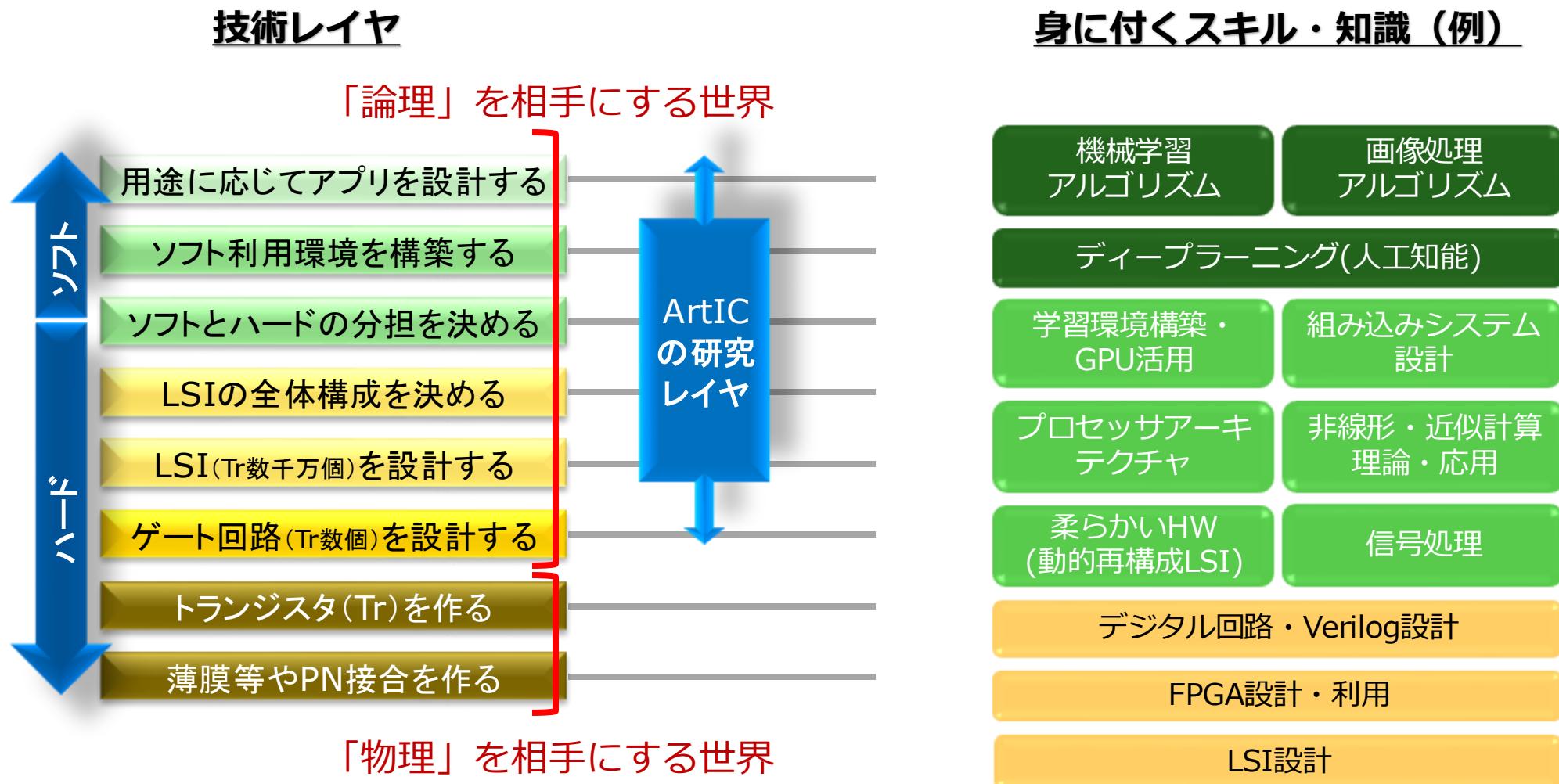
		本村研	藤木研	
教員	教授	本村 真人		全体運営、ディープラーニング
	准教授		藤木 大地	アーキテクチャ
	助教	金子 竜也		立ち上げ
	特任助教/ポスドク			
秘書		橋本, 土屋		
学生	博士	9名		
	修士	5名		
	学部	3名		

← ユニット一体運営 →

Thiem van Chu,
川村一志 のお二人が転出済み

- 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



SW-HWの両方に興味がある人向き

研究室生活

ArtICホームページ =>

- 研究に集中できるよう、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

「良い研究は良い環境から」

まだ始まつばかりの研究ユニット(研究室)です。自由で闊達な雰囲気、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、美麗で広く眺望の良いオフィス(すすかけ台・J棟17階フロア全体)という生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ティベーションが沸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 (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 处理効率化の研究、伊藤
- 亂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 (東工大 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











フロア案内

FLOOR INFORMATION

17F



1703		共通事務室1
1706		情報工学系 Department of Computer Science
1707		山村研究室 Yamamura Lab.
1710		1706 教員室(山村) Professor M. Yamamura
1707		1707 研究室
1710		1710 学生室
1704		小野 功研究室 Isao Ono Lab.
1705		1704 教員室(小野 功) Associate Professor Isao Ono
1709		1705 学生室
711		1709 研究室
711		科学技術創成研究院 Institute of Innovative Research
712		AIコンピューティング研究ユニット AI Computing Unit
713		711 ユニット控室
714		712 学生室
715		713 教員室(本村) Professor M. Motomura
716		714 教員室(劉) Associate Professor J. Yu
715		715 秘書室・会議室
716		716 教員・研究員室

← 1711-1716
1713-1714-1715

715

ARTIC



1715

ArtIC
秘書室

Professor M. Matsumura

Associate Professor J. Ito

助教等

← 1711-1716

1713-1714-1715 →

1715













ANSEL ADAMS







