

# ポストムーアの半導体技術と AIチップ設計拠点の活動

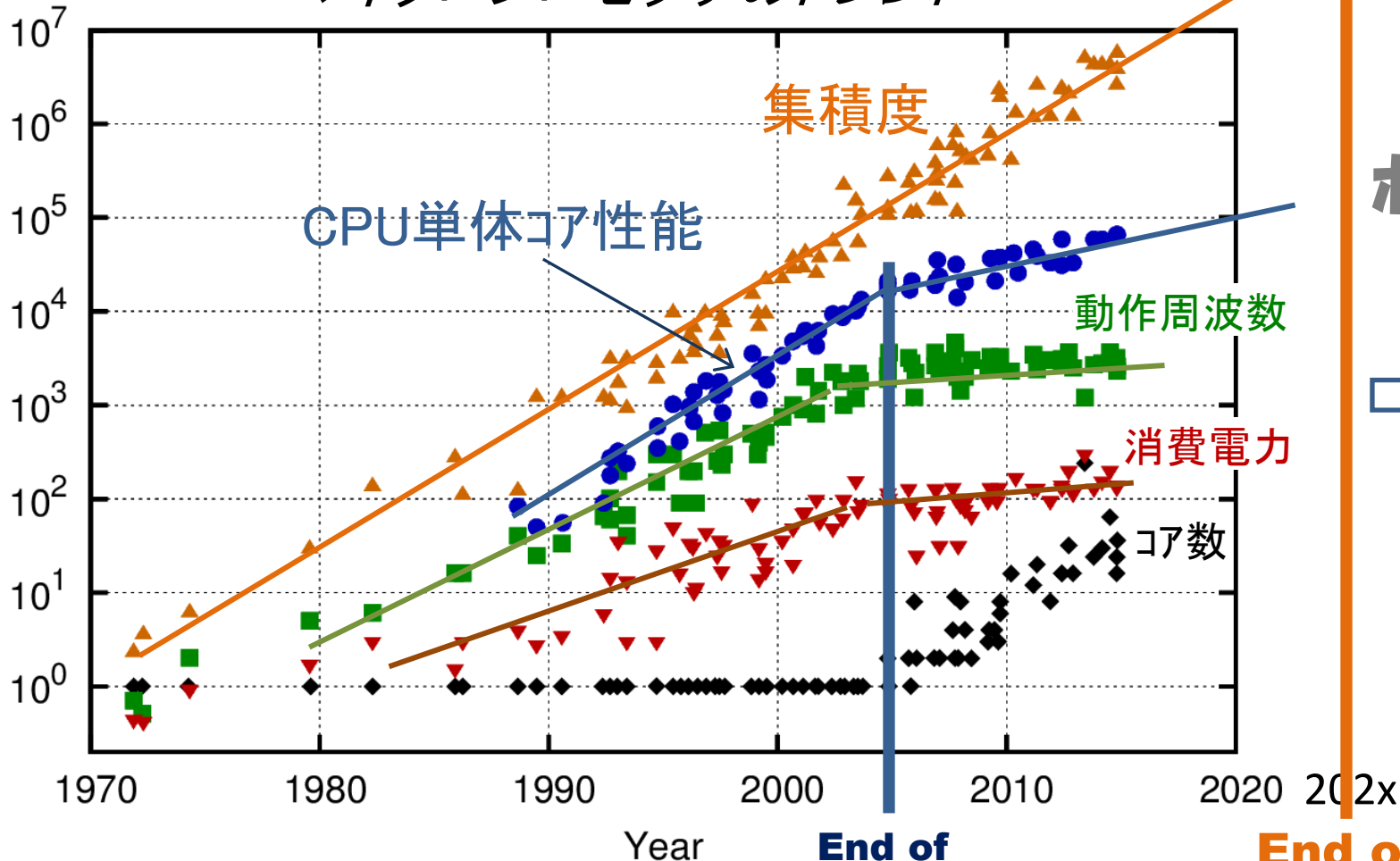
産業技術総合研究所  
内山邦男

1. ポストムーアの半導体技術

2. AIチップ設計拠点の活動

# 半導体技術の歴史的転換点

## マイクロプロセッサのトレンド



ポストムーア時代  
(202x年～)

Original data up to 2010 collected and plotted by M.Horowitz, F.Labonte, O.Shacham, K.Olukotun, L.Hammond, and C.Batten  
New plot and data collected for 2010-2015 by K.Rupp

**End of Dennard scaling**

- ・寸法:  $1/k$
- ・スイッチング速度:  $k$
- ・消費電力:  $1/k^2$

**End of Moore's Law ?**

集積度: 18カ月で2倍

# ムーアの法則の終焉後に 「指数関数的な性能向上」を担うものは何か？

- 2012 IEEE Rebooting Computing Initiativeがスタート  
“**to rethink the computer, "from soup to nuts," including all aspects from device to user interface.**”
- 2015 次世代スパコンに向けた大統領令発令  
戦略目標の1つに“**ポストムーアに向けた研究開発の強化**”を設定  
IARPA(情報機関係の研究支援組織)が量子計算、JJ計算機、  
ニューロ計算のプロジェクト推進強化
- 2016 Googleが**TPU**発表 (TPU2.0 in 2017, TPU3.0 in 2018)
- 2017 AppleがiPhone用A11チップに**Neural Engine**搭載  
DARPAが**Electronics Resurgence Initiative**プロジェクト開始
- 2018 Patterson教授のTuring Awardレクチャー@ISCA  
“**A New Golden Age for Computer Architecture**”

# International Symposium on Computer Architecture, 2017

## **In-Datcenter Performance Analysis of a Tensor Processing Unit™**

Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon,

**Many architects believe that major improvements in cost-energy-performance must now come from domain-specific hardware.**

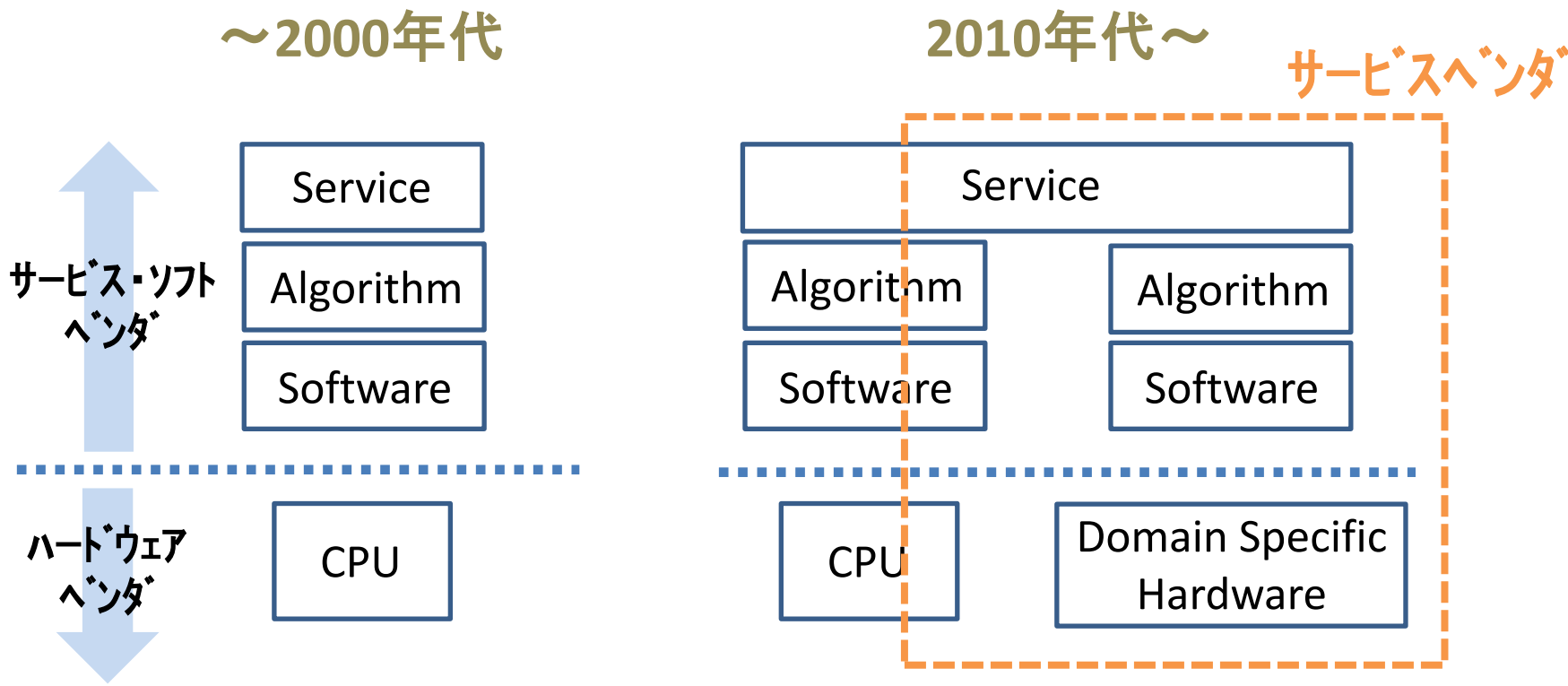
*To appear at the 44th International Symposium on Computer Architecture (ISCA), Toronto, Canada, June 26, 2017.*

### **Abstract**

**Many architects believe that major improvements in cost-energy-performance must now come from domain-specific hardware. This paper evaluates a custom ASIC—called a *Tensor Processing Unit (TPU)*— deployed in datacenters since 2015 that accelerates the inference phase of neural networks (NN). The heart of the TPU is a 65,536 8-bit MAC**

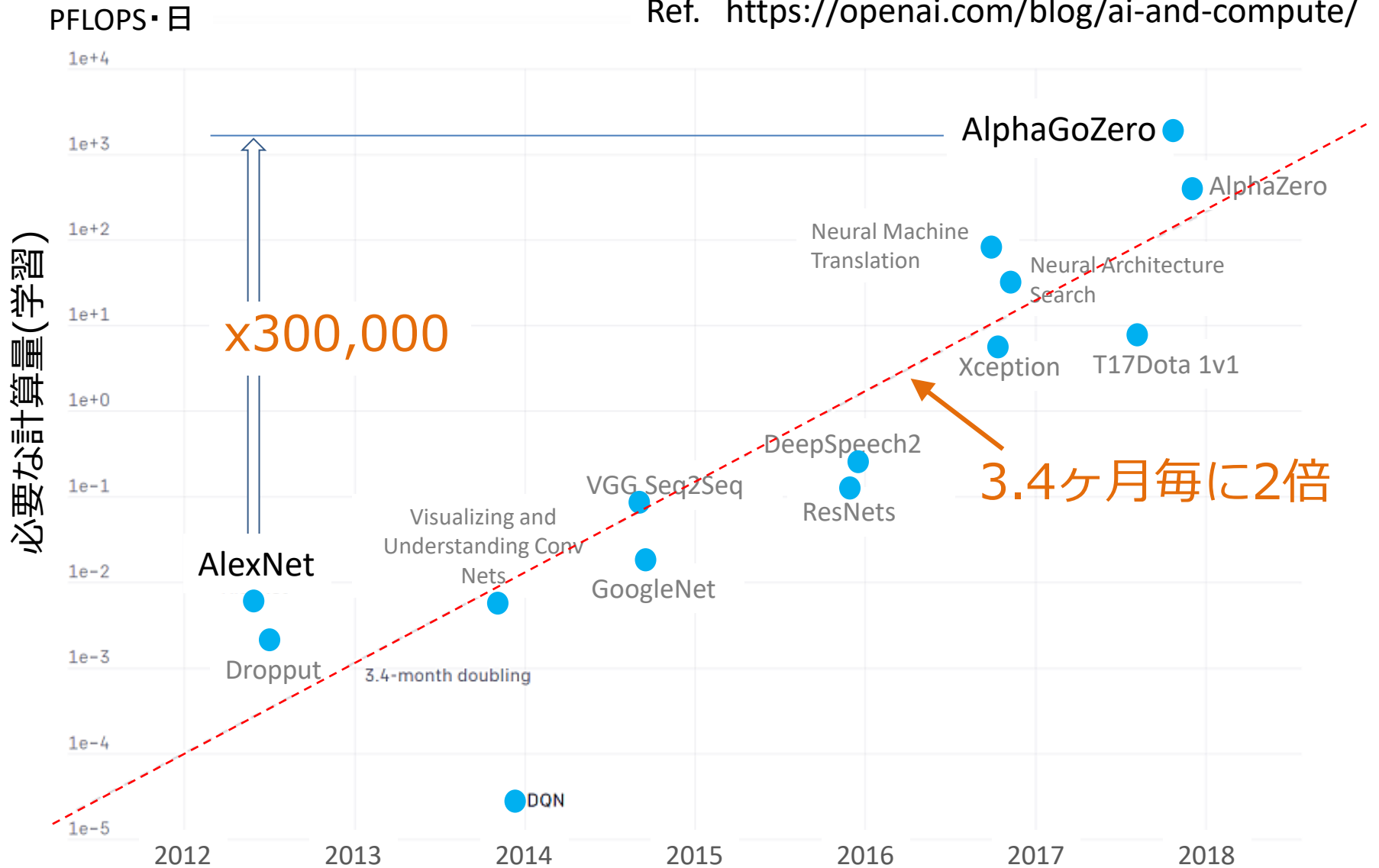
# D.Patterson教授のTuring Award受賞講演@ISCA2018

“A New Golden Age for Computer Architecture:  
**Domain-Specific Hardware/Software Co-Design**,  
Enhanced Security, Open Instruction Sets, and  
Agile Chip Development”



# AI学習( $y_i = \sum_j w_{ij}x_j$ )の計算量

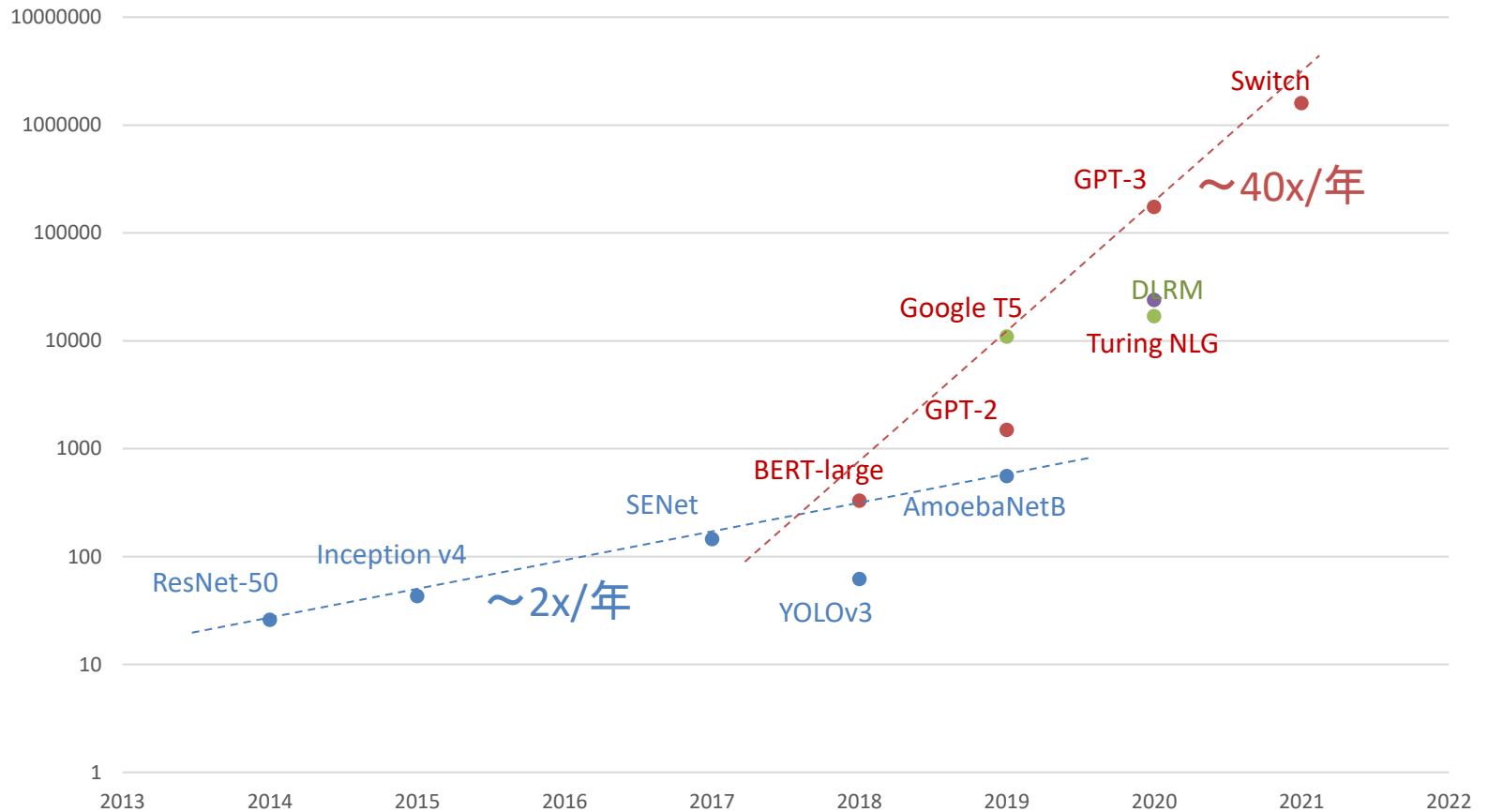
Ref. <https://openai.com/blog/ai-and-compute/>



# AI処理モデル( $y_i = \sum_j w_{ij}x_j$ )のパラメータ数

Ref. Keynote by Linley Gwennap,  
Linley Fall Processor Conference 2021

AI処理モデルパラメータ数 (million)





# サービス事業者によるAIチップ開発

- **Google**

AIチップ(TPU)を自社開発 (2015, 2017, 2018)  
自動翻訳、WEB検索などのクラウドサービスに活用

- **Apple**

iPhone用SoCを自社開発 (A4~A11, 2010~)  
A11(2017)で、AIアクセラレータを搭載

- **Facebook**

自然言語処理 (NLP) に特化したAIチップ開発中(2019/ISSCC)

- **Amazon**

推論専用AIチップ(AWS Inferentia)を開発(2018)  
450名のチップ開発チーム

- **Baidu**

クラウド向けAIチップ(Kunlun)を開発 (2017, 2020/Hot Chips)  
14nm, HBMx2, 150W, 256TOPS, PCIeGen4x8

- **Alibaba**

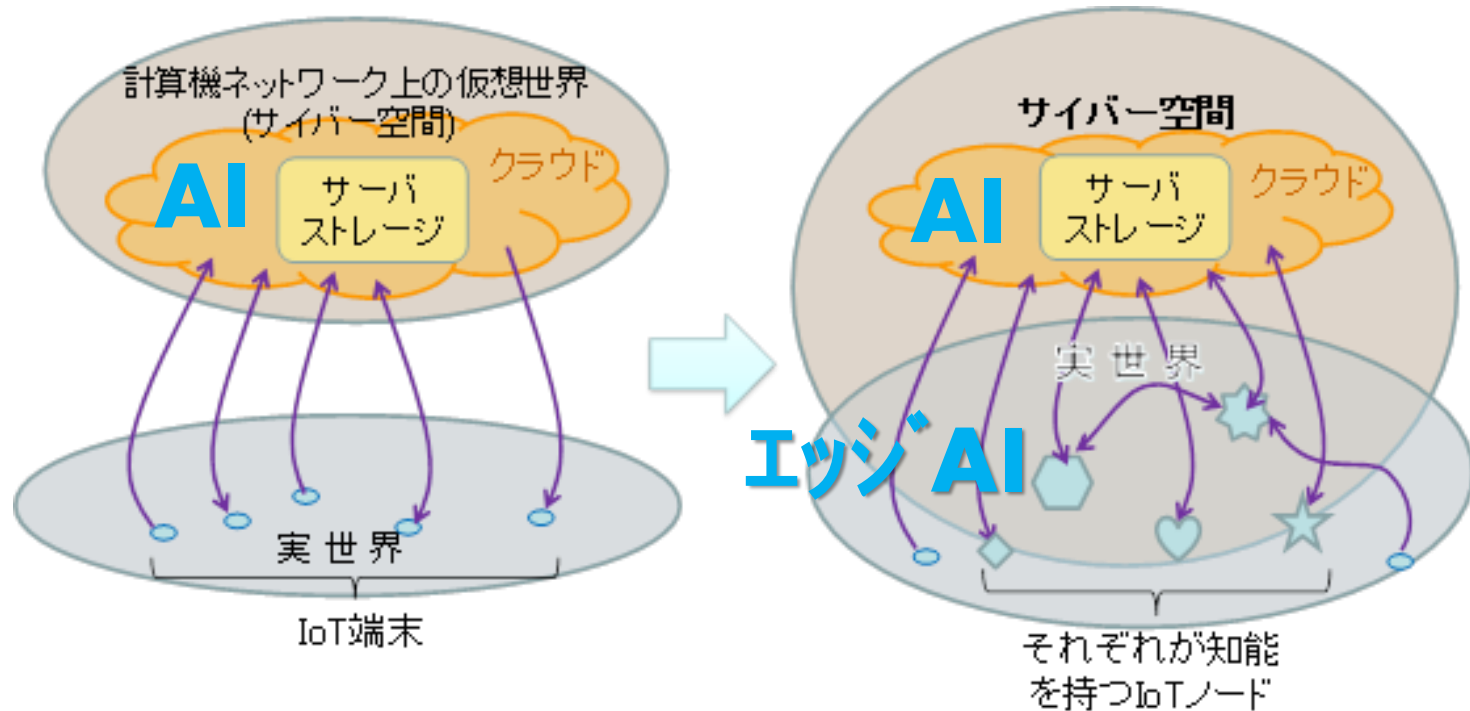
推論専用AIチップ(Hanguang 800)を開発中 (2020/Hot Chips)  
12nm, 825TOPS, 192MB on-chip SRAM, 700MHz

- **Preferred Networks**

学習用AIチップ(MN-Core)を開発 (2018, 2020/ISC)  
524TFLOPS, 500W/4die

# クラウドからエッジへ

「2030年に向けた次世代計算機技術開発戦略」 日本工学アカデミー報告書



## 「AI Moves From the Cloud to the Edge」

基調講演 by Linley Gwennap, Linley Fall Processor Conference, October 20, 2021

# エッジAIチップの状況

Ref. Linley Fall Processor Conference, 2021

- **ハイエンド (自動運転、 etc.)**

  - **Nvidia Orin X** : 72 TOPS at 25W~254 TOPS at 100W

  - **Xilinx Versal AI Edge FPGA** : 202 TOPS at 75W

  - **Qualcomm Cloud AI 100** :100 TOPS at 15W~400 TOPS at 75W

- **ミドル/ローエンド (監視カメラ、 etc.)**

  - **Ambarella** : 12 TOPS accelerator at 5W

  - **Quadric** : 4 TOPS at 2W

  - **ArchiTek** : 4 TOPS at 1.5W

  - **RealTek** : 0.7 TOPS at 0.8W

- **超低電力 (IoTセンサ、 etc.)**

  - **Syntiant** : 500uW

  - **Innatera** : 100uW

  - **Ambient, BrainChip** : サブmW

# AIチップ関連ベンチャ/SME

Horizon Robotics(中国)

Cambricon Technologies(中国)

Cerebras Systems(US)

Graphcore(UK)

Wave Computing(US)

Hailo(イスラエル)

Mythic(US)

Kneron(US)

Flex Logix Technologies(US)

DeePhi(中国)

NSITEXE(日本)

Syntiant(US)

Cornami(US)

TRIPLE-1(日本)

SiMa.ai(US)

GrAI Matter Labs(仏)

brainchip\*(US)

Eta Compute(US)

Movellus, Inc.(US)

Greenwaves Technologies(仏)

innatera Nanosystems(オランダ)

Groq(US)

edgeCortix(日本)

LeapMind(日本)

Think Silicon(ギリシャ)

Perceive(US)

Alphawave(カナダ)

Expedera(US)

Gyr Falcon Technology(US)

ArchiTeK(日本)

Digital Media Professionals(日本)

Axell(日本)

⋮

# 各国の半導体関連強化施策

第4回 半導体・デジタル産業戦略検討会議  
「半導体戦略の進捗と今後」資料より

国・地域	産業支援策の主な動向
米国	<ul style="list-style-type: none"><li>最大<b>3000億円/件</b>の補助金や「<b>多国間半導体セキュリティ基金</b>」設置等を含む国防授權法（NDAA2021）の可決。</li><li>バイデン大統領はCHIPS法案に賛意を表明。上院においては<b>5.7兆円</b>の半導体関連投資を含む「<b>米国イノベーション・競争法案</b>」が通過。</li></ul>
中国	<ul style="list-style-type: none"><li>「<b>国家集積回路産業投資基金</b>」を設置（'14, '19年）、<b>半導体関連技術へ、計5兆円を超える大規模投資</b>。</li><li>これに加えて、地方政府で<b>計5兆円を超える半導体産業向けの基金</b>が存在（<b>合計10兆円超</b>）</li></ul>
欧州	<ul style="list-style-type: none"><li>2030年に向けたデジタル戦略を発表。<b>デジタル移行（ロジック半導体、HPC・量子コンピュータ、量子通信インフラ等）に1345億€（約17.5兆円）投資等</b></li><li>製造を含む欧州の最先端チップ・エコシステムの構築を目指し、供給の安全を確保し、欧州の画期的技術のための新たな市場を発展させる「<b>新・欧州半導体法案</b>」の制定を宣言（2021.9）</li></ul>
台湾	<ul style="list-style-type: none"><li><b>台湾への投資回帰を促す補助金等の優遇策</b>を始動。ハイテク分野を中心に<b>累計で2.7兆円の投資申請</b>を受理。（2019.1）</li><li><b>半導体分野に、2021年までに計300億円の補助金</b>を投入する計画発表。（2020.7）</li></ul>
韓国	<ul style="list-style-type: none"><li><b>AI半導体技術開発への投資に1,000億円</b>を計上。（2019.12）</li><li><b>半導体を含む素材・部品・装置産業の技術開発</b>に2022年までに<b>5,000億円以上</b>を集中投資する計画を発表。（2020.7）</li><li>総合半導体大国実現のための「<b>K-半導体戦略</b>」を策定（2021.5）</li></ul>

# AIチップ / 次世代コンピューティング / 半導体関連の施策 (2018年～)

## ・ 経産省関連

AI チップ開発加速のための検証環境整備事業  
AIチップ・次世代コンピューティングの技術開発事業  
AIチップ開発加速のためのイノベーション推進事業  
ポスト5G情報通信システム基盤強化研究開発事業  
グリーンイノベーション基金事業

## ・ 文科省関連

Society5.0を支える革新的コンピューティング技術  
革新的コンピューティング技術の開拓  
次世代X-nics半導体創生拠点形成事業

## 2. AIチップ設計拠点の活動

NEDO事業

「AIチップ開発加速のためのイノベーション推進事業

研究開発項目②: AIチップ開発を加速する共通基盤技術の開発」

[https://www.nedo.go.jp/activities/ZZJP\\_100142.html](https://www.nedo.go.jp/activities/ZZJP_100142.html)

# 拠点構築の目的

- ✓ 我が国では、ベンチャー企業等を中心に、AIチップを基にした新たなビジネスを創出させる種が多数存在。
- ✓ 一方、AIチップ設計には、高額なEDAツールやIP、検証装置(エミュレータ等)が必要であり、これらがビジネス化に向けた高いハードルとなっている。
- ✓ AIチップ設計に必要な設計・検証環境を整備し、イノベーション実現のためのAIチップ開発を加速する。

## 革新的AIチップ のアイデア

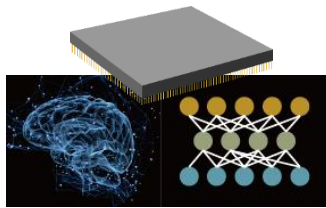


高額なEDAツール、  
IP、検証装置が必要

高いハードル

国内中小企業  
ベンチャー企業

## AIチップ プロトタイプ試作



学習、推論、認識を  
低電力かつ高速に

## 超スマート社会 (Society5.0) の実現

- ・次世代モビリティ  
自動運転, 無人配送, ...
- ・次世代ヘルスケア  
AI診断, 自動モニタリング, ...
- ・次世代サプライチェーン  
スマート保安, 無人工場, ...
- ・農林水産業スマート化  
無人農業車両, 水中ロボット, ...
- ・FinTech  
:

AIチップ設計拠点



# 拠点の体制・運営

赤字: AIチップ設計拠点 運営組織

ベンチャーキャピタル  
Pluga Capital

ベンチャー・中小企業等

EDAツール・IP・検証環境提供

サテライト拠点

ふくおかIST

**AIチップ設計拠点**

@東大本郷地区浅野キャンパス

EDAツールベンダー

IPベンダー

ファウンドリ

LSIデザインハウス  
凸版印刷, ...

ソフトウェアハウス

産総研



人工知能研究センター  
ABCI

産総研



東京大学



大学・公的機関  
コンソーシアム

- ・北海道大学
- ・東北大学
- ・福岡大学

- ✓ AIチップ設計に必要な設計・検証環境 (EDAツール, IP, エミュレータ等)の整備, 提供
- ✓ AIチップ開発に資する設計技術、検証手法の開発
- ✓ AIチップ技術に関する人材育成

# EDAツール

Cadence, Synopsys, Siemensのツールを整備

アーキテクチャ検証、高位合成、論理設計/検証、回路設計、物理設計/検証、  
論理エミュレータ、FPGAプロトタイピング、ボード設計、etc.

## IP

Synopsysの40nm, 28nm, 12nm 標準IP群

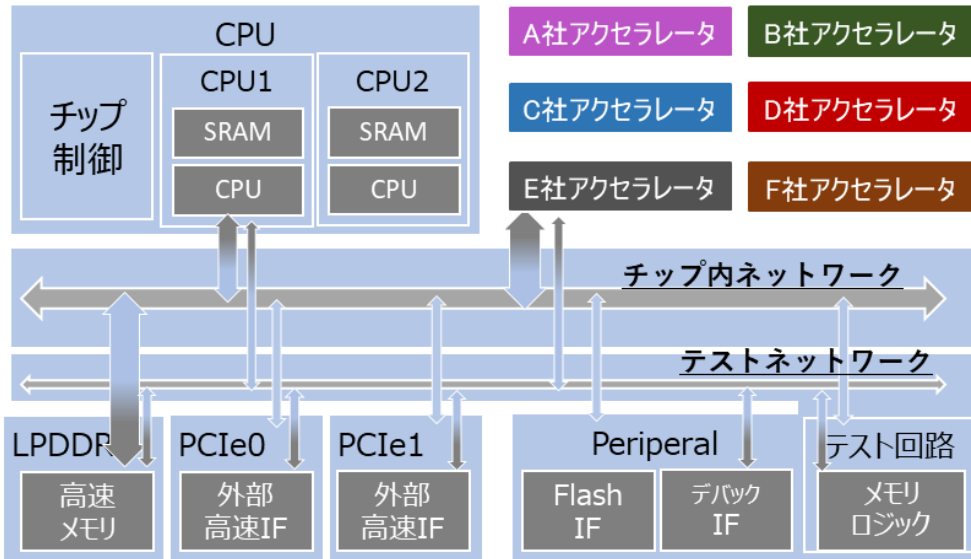
(CPU, DSP, DMAC, DDR, PCIe, USB, MIPI, I2C, etc.)

を整備

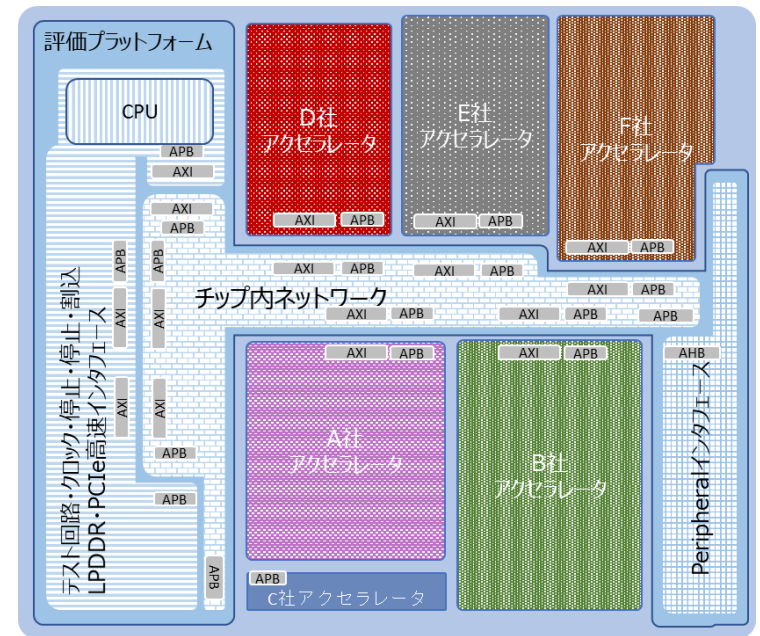
(物理系はTSMC向け)

# エッジAI向け評価チップ (AI-One)

- ・ 拠点導入IP(28nm)を活用して、拠点が評価プラットフォームを準備
- ・ 乗合チップ参加企業は各社のAIアクセラレータをプラットフォームに接続
- ・ 拠点がまとめてチップ実装を行いファブに試作依頼、各社にチップ(+ボード)を配布
- ・ 各社は試作チップ (ボード)を用いて、実証実験を実施中



チップ内部構成



チップ実装イメージ

2022年3月プレスリリース

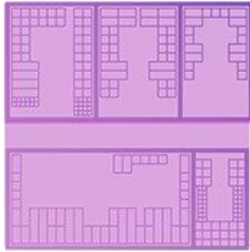
複数のAIアクセラレータを搭載した実証チップ「AI-One」の動作を確認  
-従来比45%以下の短期間で低コストのAIチップ設計・評価が可能に-

# AI-One 設計/評価プラットフォーム

## SoC設計手法

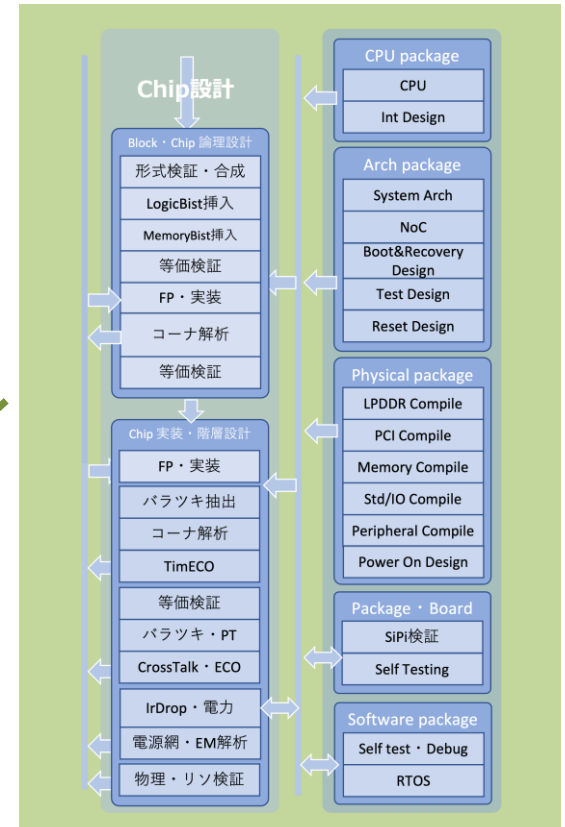
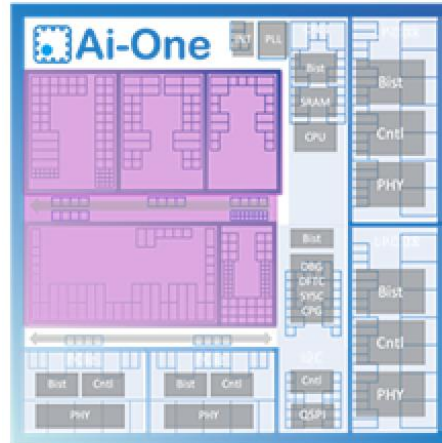
### 標準SoC回路

### AIアクセラレータ

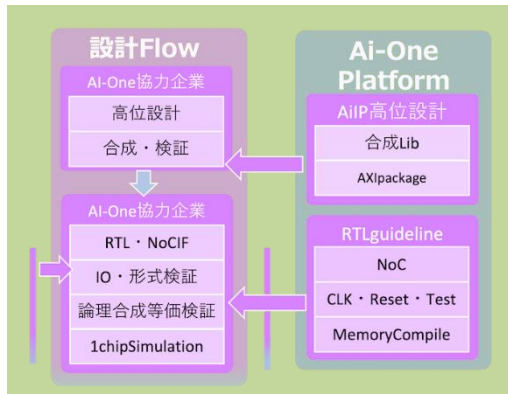


標準 SoC  
SoC化

### SoC型AIチップ



### アクセラレータ設計手法

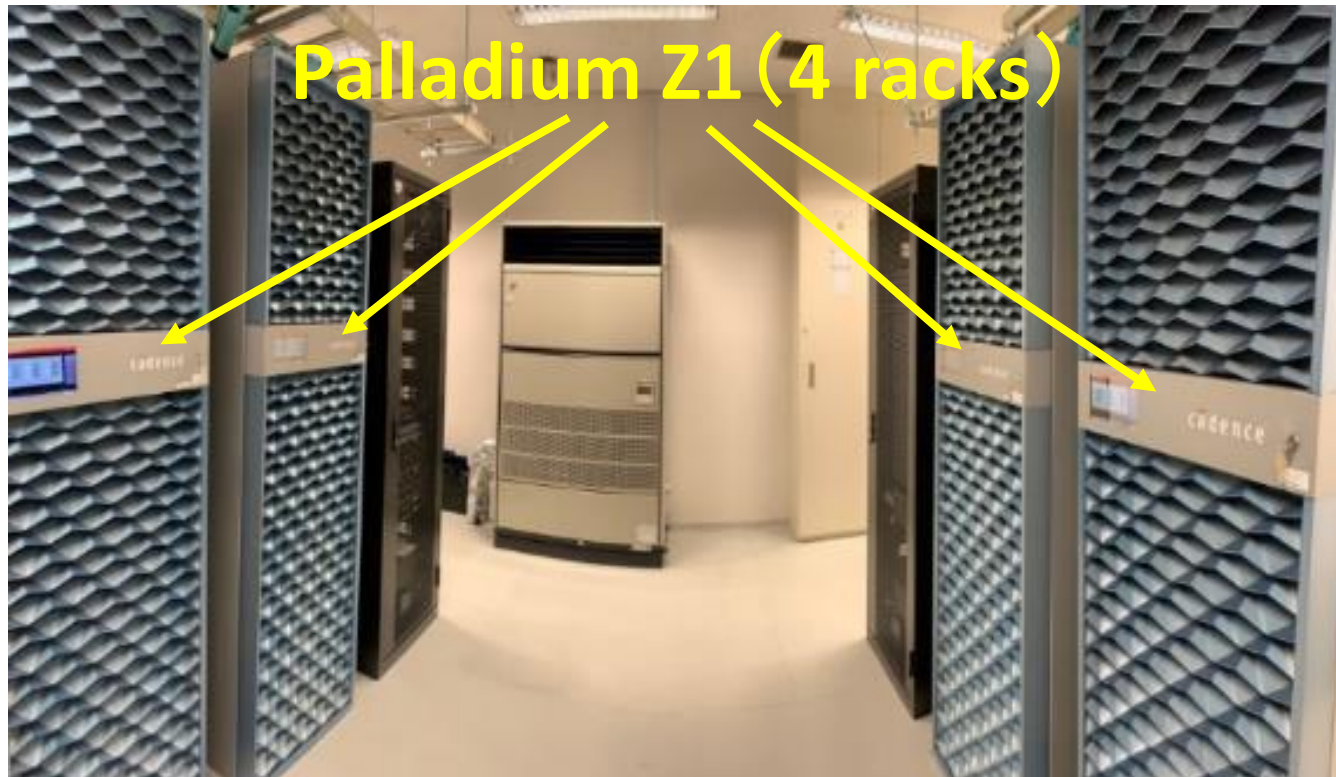


評価ボード+基盤ソフト

# 論理工ミュレータ

## Palladium Z1

- ・容量：23億ゲート, 4.6Tバイト（ユザメモリ）, 4.6Tバイト（デバッグメモリ）
- ・シミュレーション速度：最大4MHz, コパロ速度：140Mゲート/時

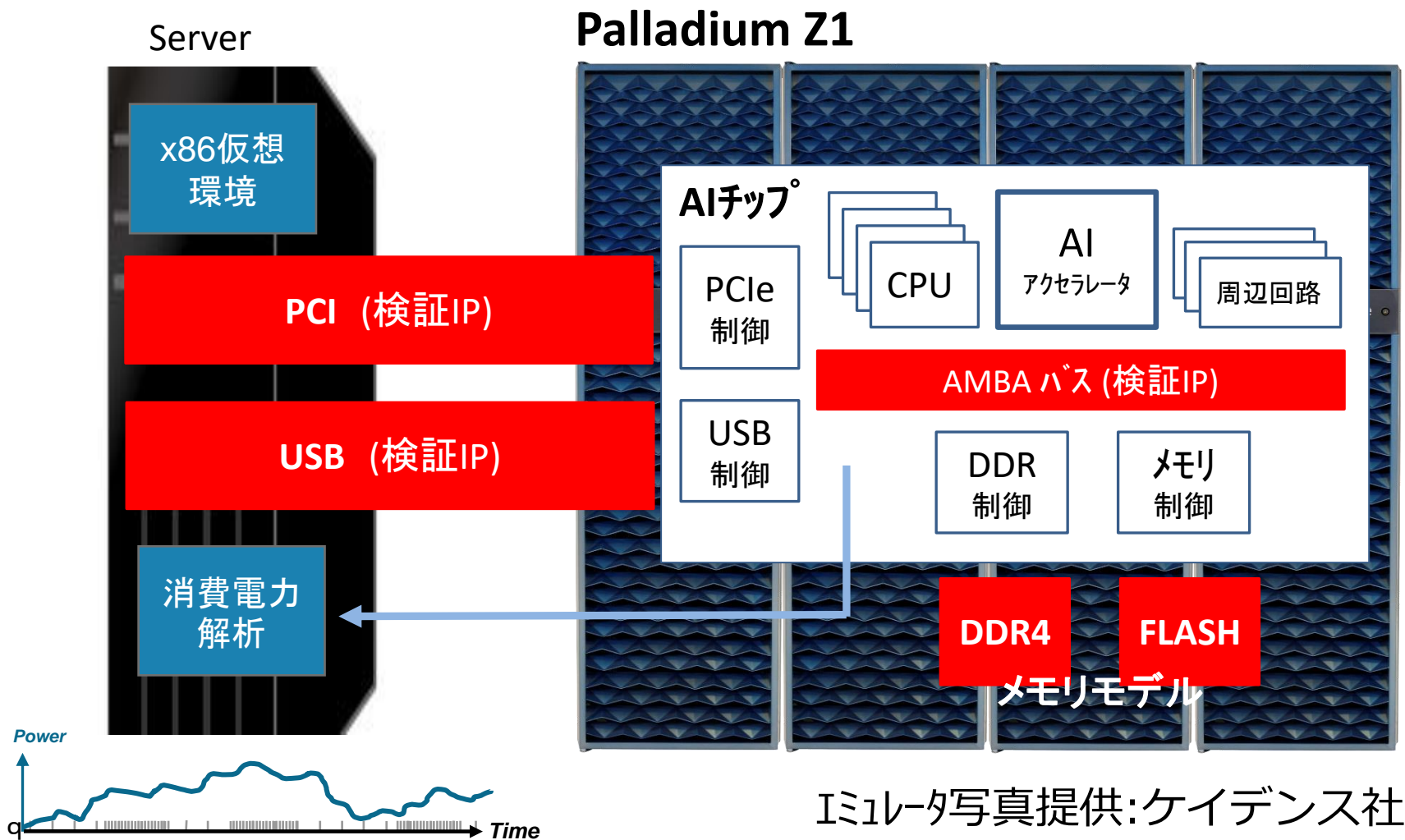


AIチップ設計拠点 サーバ室



# 論理エミュレータの活用

論理シミュレーション(ソフト)に対して4,000倍~15,000倍の高速化



エミュレータ写真提供:ケイデンス社

# RISC-V on 論理エミュレータ(Palladium Z1)

RISC-V Days Tokyo 2021 Spring, 2021/4/23, 講演

荒川 文男

東京大学 大学院 工学系研究科 システムデザイン研究センター

## 『Emulator上でのRISC-VプロセッサによるLinuxブート』

- ❑ Cadence社のEmulator上でRISC-VコアによるLinuxのブートに成功
- ❑ RISC-Vコアは東工大一色研の高位システム設計検証環境で開発したコア
- ❑ Emulation環境構築では、UARTトランザクタを活用して、Linux端末ウィンドウおよびEmulation状況のモニタリングを実現
- ❑ 約132M cycleの実行時間は最速で2分20秒程度
- ❑ 速度はCPU時間で約1.2M cycles/s、実行時間で0.9~1.1M cycles/s

# EmulationとSimulationの比較

	parameters			boot cycles		Linux Boot			1M cycles		
	ext. mem.	log check	UART			RTL	Gate	G/R	RTL	Gate	G/R
1	log	Yes	1/2	116M	Emu.	13'38"	15'06"	1.11	0'46"	2'01"	2.63
					Sim.	<i>262:46'20"</i>	<i>561:47'40"</i>	<i>2.14</i>	<i>2:15'55"</i>	<i>4:50'35"</i>	<i>2.14</i>
					S/E	<i>1,156</i>	<i>2,232</i>		177	144	
7	mem.	No	1/234	132M	Emu.	<b>2'20"</b>	<b>2'31"</b>	<b>1.08</b>	0'22"	0'53"	2.41
					Sim.	<b><i>145:09'48"</i></b>	<b><i>609:19'36"</i></b>	<b><i>4.20</i></b>	<i>1:05'59"</i>	<i>4:36'58"</i>	<i>4.20</i>
					S/E	<b><i>3,733</i></b>	<b><i>14,527</i></b>		180	314	

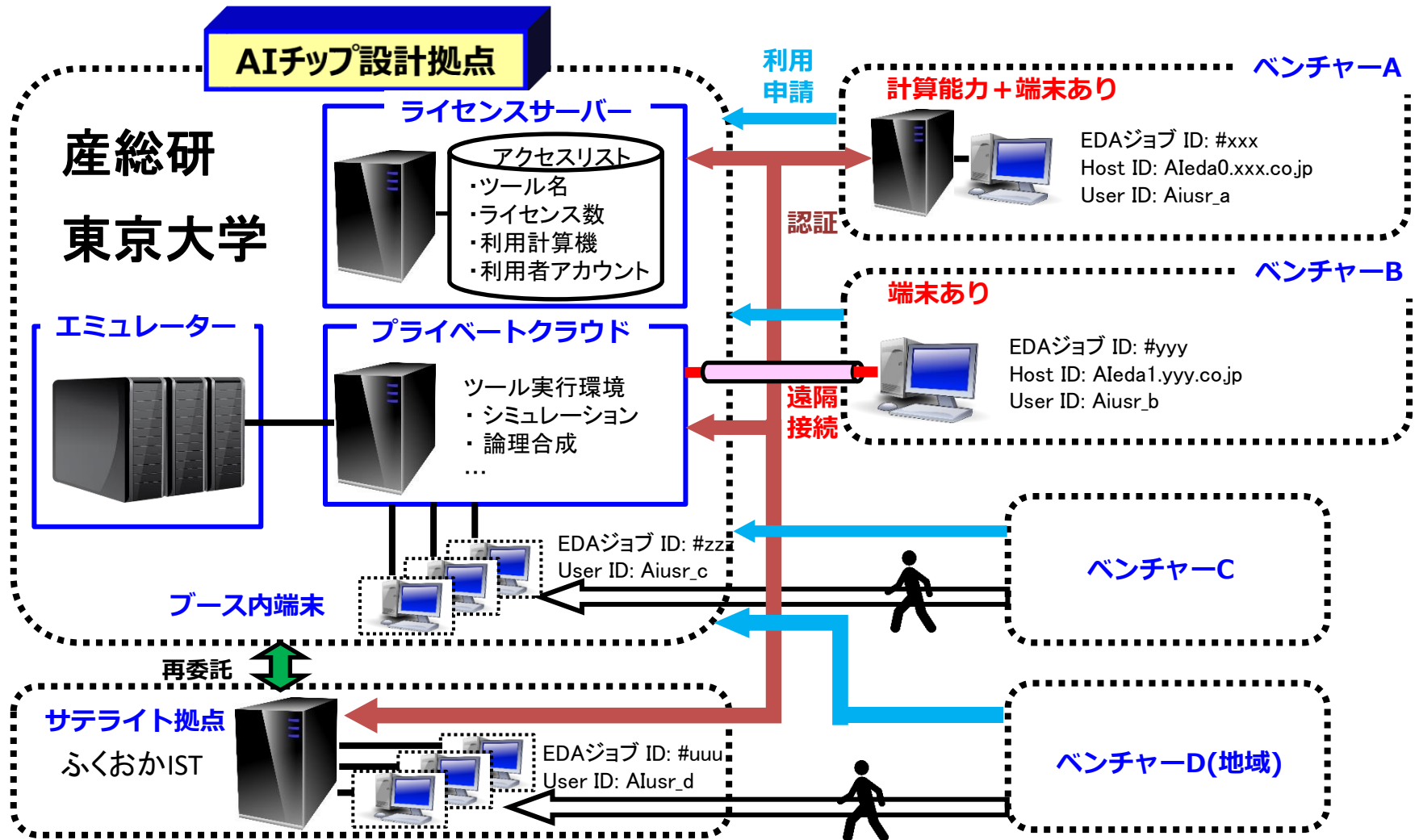
– *Italic numbers are estimated values with extrapolation of 1M-cycle executions.*

- 1M cyclesのEmulationでは初期化時間が無視できず、Linux Bootの1/100未満の実行に1/18～1/3の時間がかかる。したがって、**1M cyclesでの比較は不適切**。
- 一方、Simulationでは初期化時間が無視でき、**外挿によるLinux Boot時間推定**は正しいと期待できる。
- UART=1/2設定の影響は変化する。Emulationでは大、SimulationのRTLでは小、Gateではほぼ無し。変化の少ない**UART=1/234での比較**が妥当。
- 赤字部分がこれに相当し、**RTLで3,700倍、Gateで14,500倍高速化**。



# 拠点利用形態

- ✓ 企業毎の設計環境に応じた拠点利用形態を整備し、中小・ベンチャー等の企業群が使い易い拠点をを目指す



# AIチップ設計拠点フォーラム

- ✓ AIチップ、次世代コンピューティング、LSI設計などに関する技術情報を共有し、議論する場を提供
- ✓ 月1回のペースで開催(2019/5～)

## 第31回 AIチップ設計拠点フォーラム (2022/1/28)

- 13:30-13:35 **AIチップ設計拠点フォーラムについて**  
(産総研／内山邦男)
- 13:35-14:35 **HotChips2021にみる機械学習アクセラレータの動向**  
(東工大／本村真人先生)
- 14:35-15:35 **半導体デジタル産業戦略と関連施策について**  
(経済産業省 / 齋藤尚史氏)
- 15:40-16:40 **次世代ロジック半導体に対応したオープンイノベーション拠点の整備**  
(産総研／林喜宏氏)

# 拠点HP (https://www.ai-chip-design-center.org/)

## プロジェクトID申請には、

既にプロジェクトIDをお持ちで追加機能申請の場合は、本フォームからプロジェクトに参加される拠点利用者全員のEmailアドレスを記載願います。プロジェクトID入手後、拠点利用者全員に拠点ID申請をお願いします。(プロジェクトID申請は、管理責任者以外には必要ありません)

## プロジェクト

プロジェクトID(1)

プロジェクト名(1) **必須**

プロジェクト概要(1) **必須**

プロジェクトID(2)



AIチップ設計拠点 設計ツール Ai-One 拠点利用方法 Information ▼ Login

# イベント

拠点フォーラム  
e-講座・講演

イベント > e-講座 >

Public / お問い合わせ

お名前 **必須**

姓(例: 山田)  名(例: 太郎)

お名前(カタカナフリガナ) **必須**

姓(例: ヤマダ)  名(例: タロウ)

会社名(法人の方) **必須**

所属機関(例: 株式会社〇〇〇〇、〇〇大学)派遣元や複数所属している機関も全て記載ねがいます。

メールアドレス **必須**

例: yourname@sample.com

お問い合わせ内容 **必須**

お問い合わせ内容をご記載ください。入力は、日本語と英語の文字でお願いします。プログラムコードで使用される特殊文字を入れると問い合わせできないことがあります。

設計拠点のプライバシーポリシーに  同意する **必須**

記載内容確認。まだ送信されません。

6:40 AIチップ設計拠点フォーラム (第26回)  
2021/08/17 09:00 夏季休暇とシステム稼働  
ワーク不安定の報告(事後連絡)

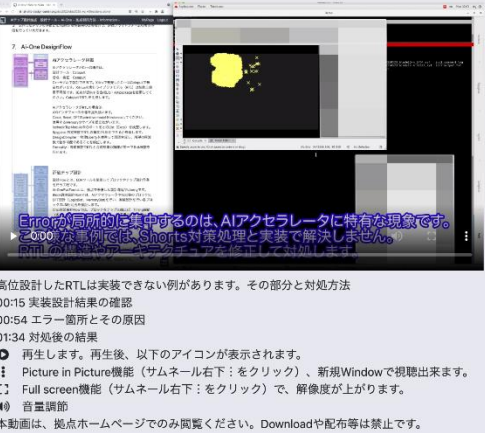
- ### 4. 教材
- Emulatorと論理Simulator, 高位設計、論理設計と検証
- エミュレータフローの一般論 (初級編)
  - デジタル設計フローの一般論 (初級編)
  - 高位合成を使ったデジタル設計 (基礎編)
  - エミュレータ論理検証の基礎 (初級編)
  - エミュレータ論理検証の基礎 (応用編)
  - A Design Verification Management Platform
  - 論理設計検証技術
  - エミュレータ活用 運営1.5時間コース
  - エミュレータ活用 関連セミナー(1)
  - エミュレータ活用 関連セミナー(2)
  - Spyglass演習1

フォーラム  
ライン  
ulation

産総研  
国立研究開発法人  
産業技術総合研究所  
エレクトロニクス・製造領域

ブースとサテライト  
東京地区・福岡地区  
FPGAにPC接続しシステムソフト実行  
設計端末with高速構内回線

## Webinarと教材Download



高位設計したRTLは実装できない例があります。その部分と対処方法

00:15 実装設計結果の確認

00:54 エラー箇所とその原因

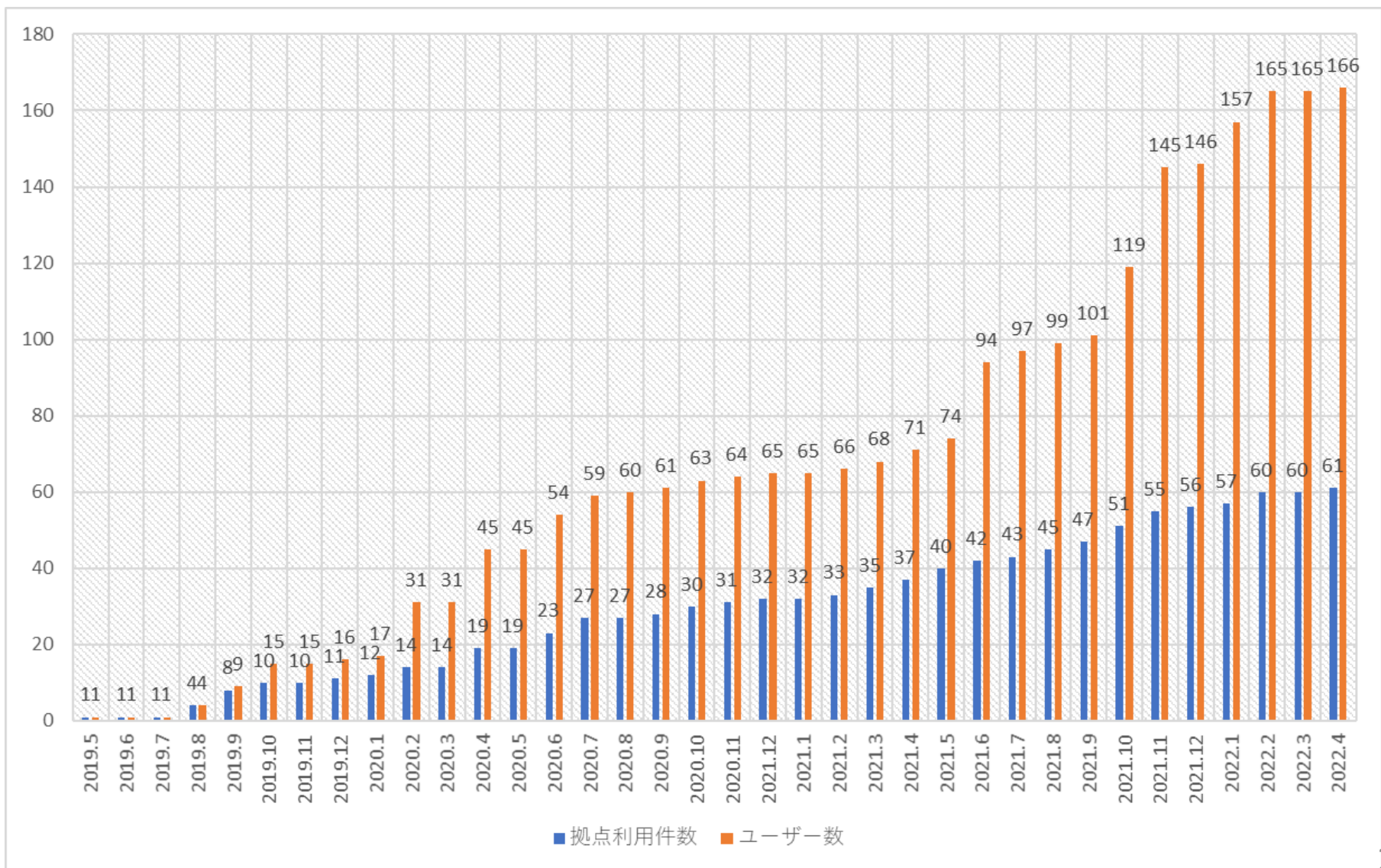
01:34 対処後の結果

- 再生します。再生後、以下のアイコンが表示されます。
- Picture in Picture機能 (サムネール右下: をクリック)、新規Windowで視聴出来ます。
- Full screen機能 (サムネール右下: をクリック) で、解像度が上がります。
- 音量調節

本動画は、拠点ホームページでのみ閲覧ください。Downloadや配布等は禁止です。

# 拠点利用件数/ユーザー数の推移

利用件数 (2022/4月時点) : 61件 (企業:43件、大学:11件、国研:6件、その他:1件)



拠点ホームページ・拠点コンタクト先

<https://www.ai-chip-design-center.org>

## AIチップ設計拠点事務局

➤ TEL: 03-5841-8460

➤ 住所: 〒113-0032

東京都文京区弥生2-11-16武田先端知ビル203号室