

コンピューティングプラットフォームの展望

- 東洋大学情報連携学部 (INIAD)
- 理化学研究所計算科学研究センター
- 清水 徹

あらまし

1. 今、コンピュータのR&Dは
2. コンピュータアーキテクチャの原点
3. コンピュータアーキテクチャの実装
4. 高性能化のチャレンジ
5. 何のための高性能化か
6. コンピューティング変革の兆し
7. コンピューティング周りのトレンド
8. まとめ

1. 今、コンピュータのR&Dは

プロセッサのR&D動向

● Hot Chips 2020, Day 1

■ サーバプロセッサ

- Next Generation Intel Xeon(R) Scalable Server Processor: Icelake-SP
- IBM's POWER10 Processor
- Marvell ThunderX3™ Next Generation Arm-Based Server Processor
- The 5.2GHz IBM z15 Processor

■ モバイルプロセッサ

- AMD Next Generation 7nm Ryzen(TM) 4000 APU
- Inside Tiger Lake: Intel's Next Generation Mobile Client CPU

プロセッサのR&D動向

● Hot Chips 2020, Day 1

■ エッジコンピューティングとセンシング

- Xuantie-910: Innovating Cloud and Edge Computing by RISC-V, Alibaba
- A technical overview of the Arm Cortex-M55 and Ethos-U55: ARM's most capable processors for endpoint AI
- PGMA: A Scalable Bayesian Inference Accelerator for Unsupervised Learning, Harvard University

■ GPUとゲーミングアーキテクチャ

- NVIDIA's A100 GPU: Performance and Innovation for GPU Computing
- The Xe GPU Architecture, Intel
- Xbox Series X System Architecture, Microsoft

プロセッサのR&D動向

● Hot Chips 2020, Day 2

■ FPGAとリコンフィギュラブルアーキテクチャ

- Agilix Generation of Intel FPGAs
- Xilinx Versal Premium Series
- Compute Substrate for Software 2.0, Tenstorrent

■ ネットワーキングと分散システム

- Tofino2 - A 12.9Tbps Programmable Ethernet Switch, Intel/Barefoot
- Pensando Distributed Services Architecture, Pensando
- The DPU: A New Category of Microprocessor, Fungible
- High-density Multi-tenant Bare-metal Cloud with Memory Expansion SoC and Power Management, Alibaba

プロセッサのR&D動向

● Hot Chips 2020, Day 2

■ 機械学習(ML)トレーニング

- Google's Training Chips Revealed: TPUv2 and TPUv3
- The Second Generation Cerebras Wafer Scale Engine, Cerebras
- Manticore: A 4096-core RISC-V Chiplet Architecture for Ultra-efficient Floating-point Computing, ETH Zurich

■ 機械学習(ML)推論

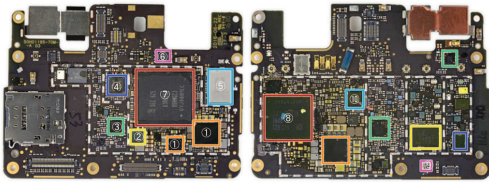
- Baidu Kunlun - An AI Processor for Diversified Workloads, Baidu
- Hanguang 800 NPU - The Ultimate AI Inference Solution for Data Centers, Alibaba
- Silicon Photonics for Artificial Intelligence Acceleration, Lightmatter

2. コンピュータアーキテクチャの原点

コンピュータアーキテクチャ

● 何のボードでしょう？

東洋大学INIAD (情報連携学部) 講義「コンピュータアーキテクチャ」から



- ① Power management IC Quick charge IC
- ② Smart audio amplifier
- ③ LTE RF transceiver
- ④ NFC controller
- ⑤ Wi-Fi
- ⑥ Low power IMU
- ⑦ Processor & mobile DRAM
- ⑧ 32 GB Universal Flash Storage
- ⑨ Pressure sensor
- ⑩ Audio codec

コンピュータアーキテクチャ

● 富士通リレー式コンピュータ FACOM128B



富士通沼津工場池田敏夫記念室にて - 富士通DNA館のご好意により

コンピュータアーキテクチャ

● Google Phone: Pixel XL Board

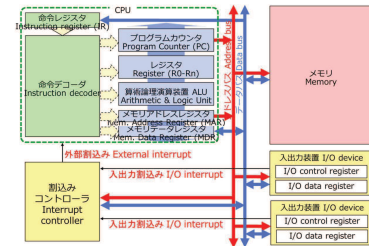
<https://www.ifixit.com/Teardown/Google+Pixel+XL+Teardown/71237>



コンピュータアーキテクチャのデザイン

● 現在のすべてのコンピュータの基本構造

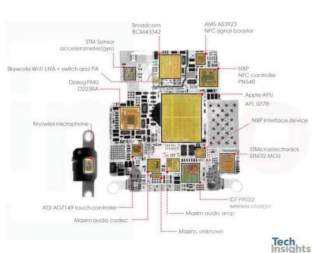
■ フォンノイマン・コンピュータ



コンピュータアーキテクチャ

● 何のボードでしょう？

東洋大学INIAD (情報連携学部) 講義「コンピュータアーキテクチャ」から



コンピュータアーキテクチャのデザイン

● フォンノイマン・コンピュータの基本

■ CPUとメモリが1本のバスで接続

- 命令やデータの位置は、CPUが生成するアドレスで指定される
 - **メモリには命令とデータが格納**されバスを介して読み書きされる
 - CPUは、メモリに格納された命令コードをデータとして読み書きできる
 - **プログラムでプログラムを編集、変換、合成**することができる
 - 従って、コンピュータの機能をソフトウェアで自己拡張できる
- ➔ **プログラミングモデル**を定めている

コンピュータアーキテクチャ

● Apple watch S1 processor

<http://www.techinsights.com/abouttechinsights/overview/blog/apple-watch-teardown/>



3. コンピュータアーキテクチャの実装

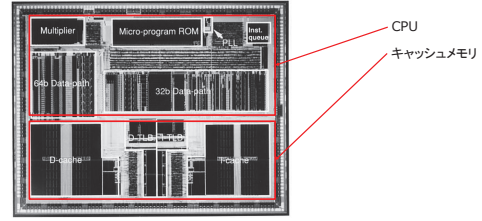
コンピュータアーキテクチャの実装

- フォンノイマン・コンピュータ開発の方針
 - プロセッサ (CPU) とメモリが分離、バスで接続
 - プロセッサとメモリ、それぞれ専門にR&D
 - LSIチップの高集積化 (Moore's law) でR&Dを牽引
- チップ高集積化 (Moore's law) の効果
 - 回路規模 → チップ面積 → 製造コスト
 - 回路遅延 → クロック周期 → 動作周波数
 - スwitching回数×トランジスタ数 → 消費電力

17

マイクロプロセッサの30年

- 1993: Main-frame CPU "PXB1" 800nm 40MHz
 - 40MHz以上ではキャッシュは不可欠 Cache is essential above 40MHz.

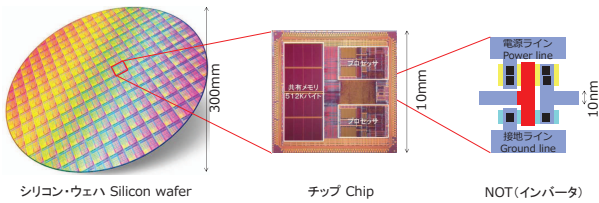


"A 1.71 M-transistor CMOS CPU chip with a testable cache architecture", IEEE ISSCC (1993)

21

コンピュータアーキテクチャの実装

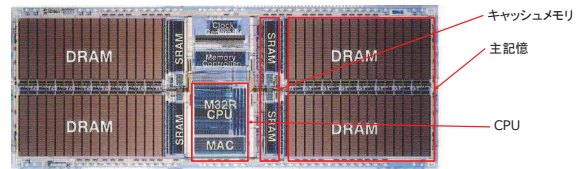
- ウェハ Wafer → チップ Chip → ゲート Gates



18

マイクロプロセッサの30年

- 1996: "M32R/D" 450nm 66MHz
 - DRAMのメインメモリ内蔵 DRAM main memory is embedded.

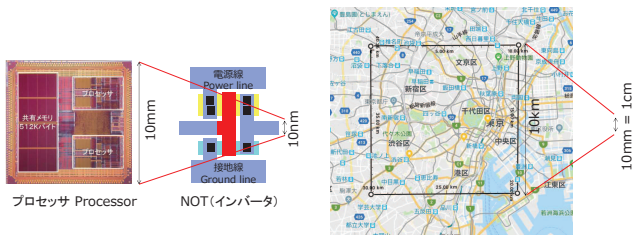


"A multimedia 32 b RISC microprocessor with 16 Mb DRAM", ISSCC (1996)

22

コンピュータアーキテクチャの実装

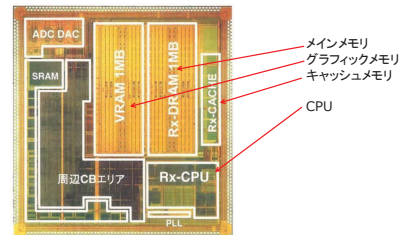
- LSIチップ設計の規模と精度



19

マイクロプロセッサの30年

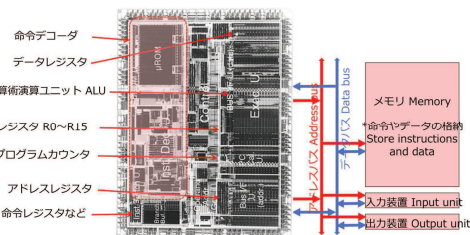
- 1999: "A camera chip with M32R/D" 180nm 90MHz
 - カメラ画像処理内蔵 Built-in camera image processing.



23

マイクロプロセッサの30年

- 1990: TRON "Gmicro/100" 1000nm (1μm) 25MHz

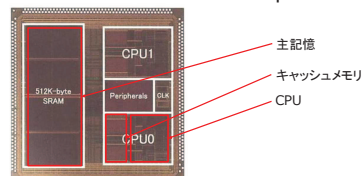


"The Gmicro/100 32-Bit Microprocessor", IEEE Micro (Aug. 1991)
「32ビットMPU Gmicro/100」, 日経エレクトロニクス (1989.7.10)

20

マイクロプロセッサの30年

- 2003: "M32R multicore" 150nm 600MHz
 - 対称型のマルチコアを内蔵 Symmetric multicore in a chip.
 - Linux SMPを移植 Linux SMP is ported.



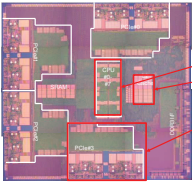
"A 600 MHz single-chip multiprocessor with 4.8 GB/s internal shared pipelined bus and 512 kB internal memory", ISSCC (2003)

24

マイクロプロセッサの30年

● 2011: "M32R multicore" 45nm 400MHz

- マルチポートのネットワーク制御用チップ
Multiport network control chip



(CPU+1次キャッシュメモリ) x 8コア
2次キャッシュメモリ
高速バスのインタフェース x 4

"An 80Gb/s Dependable Communication SoC with PCI Express I/F and 8 CPUs", ISSCC (2011)

マイクロプロセッサの30年

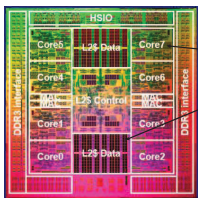
● マイクロプロセッサ設計の方針

- 微細化を「高性能×低電力」より「必要性能×多数×超低電力」に
- 「数」を「質」「精度」「機能」に
 - ソフトウェア化
- ネットワーク階層化によるデータトラフィック低電力化
 - フレキシビリティ (汎用性) が課題

マイクロプロセッサの30年

● 2010: Fujitsu "SPARC64 VIIIfx" 45nm 2000MHz

- High performance computer (HPC) 京 K's CPU chip



(CPU+1次キャッシュメモリ) x 8コア
2次キャッシュメモリ: 6M-byte

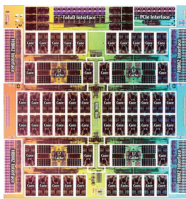
"Sparc64 VIIIfx: A New-Generation Octocore Processor for Petascale Computing", IEEE Micro (2010)

4. 高性能化のチャレンジ

マイクロプロセッサの30年

● 2018: Fujitsu "A64FX", 7nm

- HPC 富岳 Fugaku's CPU chip

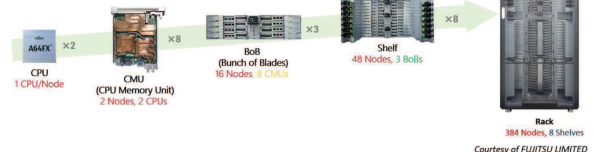


"Fujitsu High Performance CPU for the Post-K Computer", Hot Chips (2019)
「スーパーコンピュータ富岳の開発」, サイエンスフィク・システム研究会HPCフォーラム (2019)

高性能化のチャレンジ

● 富岳の高性能化技術

- 1ノードあたり、1CPU (52コア)
- 1ラックあたり、384ノード (19,968コア)
- 富岳は15万ノード超 (2,995,200,000コア)

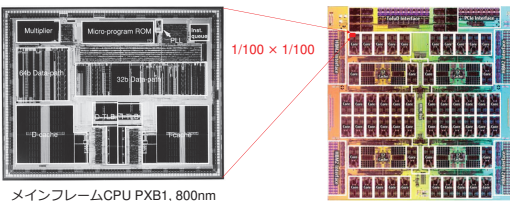


Rack
384 Nodes, 8 Shelves
Courtesy of FUJITSU LIMITED

マイクロプロセッサの30年

● マイクロプロセッサのこれから

- CPUを「ネジ、釘」として
- フレキシブルな階層コンピューティング



メインフレームCPU PXB1, 800nm

スパコンCPU A64FX, 7nm

高性能化のチャレンジ

● 高性能化の基本戦略

- 動作周波数の向上
 - チップ製造技術の微細化
 - パイプライン処理
- 並列処理
 - 命令並列
 - 演算並列
 - マルチスレッド (スレッド並列)
 - マルチコア (CPUが複数)
 - マルチコンピュータ (コンピュータが複数)
- プロセッサ-メモリ間バスの高性能化
 - 内蔵キャッシュメモリの大容量化
 - 高周波数化、多ピン化

高性能化のチャレンジ

● 高性能化の基本戦略

■ 動作周波数の向上

- チップ製造技術の微細化 → 全てのプロセッサ
- パイプライン処理 → 全てのプロセッサ

■ 並列処理

- 命令並列 → 全てのプロセッサ
- 演算並列 → GPU、DSP (Digital Signal Processor)、スパコン
- マルチスレッド (スレッド並列) → 今や、全てのプロセッサ
- マルチコア (CPUが複数) → 今や、全てのプロセッサ
- マルチコンピュータ (コンピュータが複数) → クラウド・サーバー、スパコン

■ プロセッサ-メモリ間バスの高性能化

- 内蔵キャッシュメモリの大容量化 → 全てのプロセッサ
- 高周波数化、多ピン化 → 全てのプロセッサ

33

5. 何のための高性能化か

37

高性能化のチャレンジ

● 高性能化の壁：消費電力と集積度（コスト）



34

何のための高性能化か

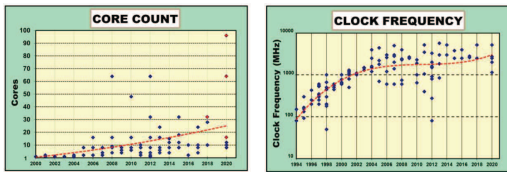
- コンピュータの高性能化は必ず新しいアプリケーションを生んできた
- 従来の方針で高性能化が難しくなってきた今、何のための高性能化が再確認することは重要
- 超高性能コンピューティングはなぜ必要か？

38

高性能化のチャレンジ

● 高性能化の壁：消費電力と集積度（コスト）

■ コア数と動作周波数の傾向



ISSCC 2020 Press kit より

39

何のための高性能化か

- “The Age of Robots”, Hans P. Moravec, CMU (1993) より
 - 2000-2010: 1,000 MIPS (第1世代汎用ロボット)
 - 爬虫類クラス … トカゲぐらい？
 - 汎用の知覚、操作、機動性 (General-purpose perception, manipulation and mobility)
 - 2010-2020: 30,000 MIPS (第2世代汎用ロボット)
 - 哺乳類クラス … ネズミぐらい？
 - 適応性が高い学習 (Accommodation learning)
 - 2020-2030: 1,000,000 MIPS (第3世代汎用ロボット)
 - 霊長類クラス … サルぐらい？
 - フールドモデリング (World modeling)
 - 2030-2040: 30,000,000 MIPS (第4世代汎用ロボット)
 - 人間クラス
 - 解析・推論・理論化 (Reasoning)

39

高性能化のチャレンジ

● 高性能化の「壁」

■ 動作周波数の向上

- LSIチップ製造技術の微細化 → 消費電力増大、チップコスト高騰
- 命令パイプライン処理 → 多段化に限界 (パイプラインハザード)

■ 並列処理

- 命令並列 → 並列化に限界 (パイプラインハザード)
- 演算並列 → 並列プログラミング
- マルチスレッド (スレッド並列、タスク並列) → 並列プログラミング
- マルチコア (CPUが複数) → 並列プログラミング、並列処理OS
- マルチコンピュータ (コンピュータが複数) → 並列処理OS

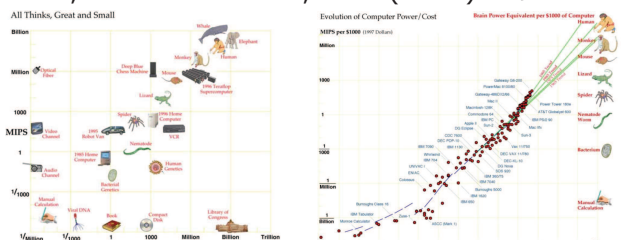
■ プロセッサ-メモリ間バスの高性能化

- 高周波数化、多ピン化 → 高周波数・多ピン実装

36

何のための高性能化か

- “When will computer hardware match the human brain?”, Hans P. Moravec, CMU (1997) より



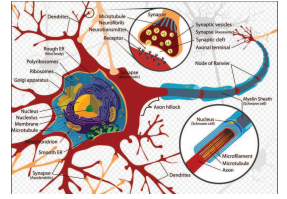
40

6. コンピューティング変革の兆し

41

コンピューティング変革の兆し

- ニューラルネットワーク
 - ニューロンとその接続（シナプス）でコンピューティング
 - メモリ：ニューロンがシナプスに信号を伝える重み
 - プロセッサ：シナプスからの刺激でニューロンが発火
 - ネットワーク：ニューロン間の接続関係
 - プログラミングが課題
 - プログラミング=トレーニング
 - プログラム=トレーニングデータ
 - 実行=推論
 - 実現方法が課題
 - デジタル：GPU, ASIC, CIM?
 - アナログ：CIM, Spiking NN?



コンピューティング変革の兆し

- 命令・プログラムの並列処理、高速化は限界
 - 命令並列
 - 演算並列
 - マルチスレッド（スレッド並列）
 - マルチコア（複数のCPU）
 - マルチコンピュータ（複数のコンピュータ）
- 更に細粒度の並列処理をコンピューティングに活用
 - リンコンフィギュラブル・コンピューティング
 - ニューラルネットワーク・コンピューティング
 - 量子コンピューティング

42

コンピューティング変革の兆し

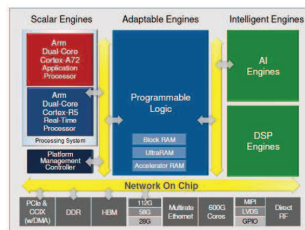
- 量子コンピュータ、イジングマシン
 - 組合せ最適化問題の並列処理ソリューション
 - 量子による高度な並列処理
 - 量子ビット数増の段階
 - 電力や集積度は議論以前
 - 汎用性はどの程度か？



日経 2019.11.23より

コンピューティング変革の兆し

- リンコンフィギュラブル・プロセッサ
 - プログラムを直接、論理回路にマッピング
 - データフローには適
 - コントロールフローには課題
 - 電力と集積度は課題
 - プログラミング・モデルは？



"Xilinx First 7nm Device: Versal AI Core (VC1902)", Hot Chips 2019より

43

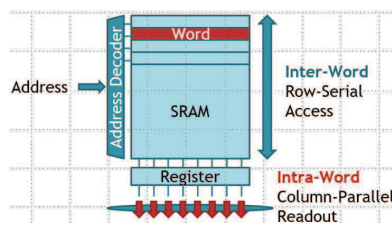
コンピューティング変革の兆し

- 細粒度の並列処理のプログラミング
 - リンコンフィギュラブル・コンピューティング
 - ニューラルネットワーク・コンピューティング
 - 量子コンピューティング
- プログラミングモデルの変革は避けられない
 - 小変革？
 - 大変革？
 - ミックス？

44

コンピューティング変革の兆し

- コンピューティング・インメモリ（CIM）
 - ワード内セルの並列処理
 - デジタルもアナログもある
 - アプリ特化：
 - DNN (Deep Neural Network)
 - アンニリング
 - メモリ技術の継承・拡張
 - デジタル：コストが課題
 - アナログ：安定性が課題



参照：2019 Riken International Workshop on Neuromorphic Computing
"AI on the Edge: Frontiers for Energy-Efficient Hardware Architectures", 本村先生より

45

7. コンピューティング周りのトレンド

プラットフォームのオープン化

- RISC-V
- プラットフォームのソフトウェア技術はオープン化してきた

Open Interfaces Work for Software!

Field	Open Standard	Free, Open Implement.	Proprietary Implement.
Networking	Ethernet, TCP/IP	Many	Many
OS	Posix	Linux, FreeBSD	M/S Windows
Compilers	C	gcc, LLVM	Intel icc, ARMcc
Databases	SQL	MySQL, PostgreSQL	Oracle 12C, M/S DB2
Graphics	OpenGL	Mesa3D	M/S DirectX
ISA	??????	-----	x86, ARM, IBM360

Why not successful free & open standards and free & open implementations, like other fields?

アプリの時代

- IoT - 森川研・川原研フォーラムより (2019)
- IoTは、ネットワーク、給電、センシング、アプリのアイデア

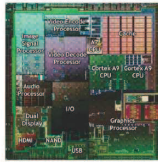
M-01 CT4oRa: Enabling High-efficient and Omnipresent IoT Network	U-01 Multimoda QSCR: 屋外中どこでもワイヤレス充電
M-02 LoRaを用いた位置情報に基づく屋内外位置測定	U-02 Atvax: 給電領域を自由に構成可能な2次元無線給電システム
M-03 電波環境マップ構築に向けた空間補間手法の性能改善	U-03 切替により野炊の良さが可能な無線電力伝送シート
M-04 遠距離LEDベーススマートフォンのメタ通信	U-04 REVOLUTIONARY 200X! : Speculative Design Toolkit for Beginners
M-05 協調機械学習におけるオンラインベース処理の軽量化	U-05 温度差発電駆動型土壌センサーの開発
M-06 降雨と気象観測データの予測性に基づく事前危険解析を用いた早期災害予測	U-06 poimo - Portable and Inflatible Mobility -
M-07 SCADA データを用いた風車異常検知	U-07 パーツナルモビリティ向け無線給電ステーション
M-08 直感音声特徴量ベースの分離再構成音声品質一	U-08 形状記憶合金で駆動するイモシム型ソフトロボット
M-09 海外調査を想定した同時実行システムの一括分岐Q学習エージェントによる学習	U-09 ヘビの鱗から着想を得た摩耗異方性表皮を持つソフト・ロボット
M-10 音楽活動における意思決定システムの開発および応用	U-10 銀ナノインク印刷で作る折り紙スピーカー
M-11 Liquid Holdup Prediction in Oil and Gas Wells	U-11 センサーとアクチュエータの一括印刷で作るソフトロボット
Using a Decision Tree Regression Predictive Model	U-12 鋼ケーブル論理回路による折り紙ロボットの制御
M-12 ナップレス RFID への情報書き込み技術	U-13 ソフトアクチュエータの音響特性を利用した長きセンシング
M-13 漸次離層型ワイヤレス給電	
M-14 メッシュ帯域近距離解析によるワイヤレス給電装置の金属筐体自動設計	

プラットフォームのオープン化

- RISC-V
- 1つのSoCにたくさんの命令セットアーキテクチャ (ISA)

Today, many ISAs on one SoC

- Applications processor (usually ARM)
- Graphics processors
- Image processors
- Radio DSPs
- Audio DSPs
- Security processors
- Power-management processor
- > dozen ISAs on some SoCs - each with unique software stack



NVIDIA Tegra SoC

Why?
 • Apps processor ISA too big, inflexible for accelerators
 • IP bought from different places, each proprietary ISA
 • Engineers build home-grown ISA cores

アプリの時代

- eスポーツ - 横浜ITクラスター交流会より
- eスポーツ = ゲームのスポーツ化、スポーツのゲーム化



プラットフォームのオープン化

- RISC-V
- RISC-Vは、オープンソースのISA
- なぜISAをオープンにするか
 - ソフトウェア開発・蓄積の効率化
 - ライセンス・フリー
 - オーナー企業依存のリスク削減
 - 品質向上のスピードアップ
 - セキュリティ向上
 - ツールなどのサポートコスト削減

アプリの時代

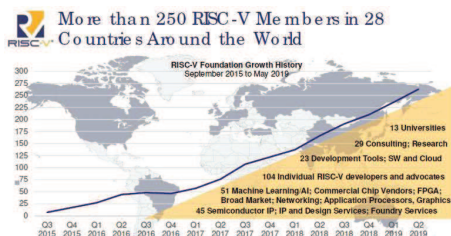
- eスポーツ
- それを実現するためのセンサー、IoT

力と動きをリアルタイムに可視化するデバイス

- 筋電位センサ: 筋肉の発火を測る
- 9軸センサ: 動き、回転、姿勢を検知する
- Bluetooth: 無線通信

プラットフォームのオープン化

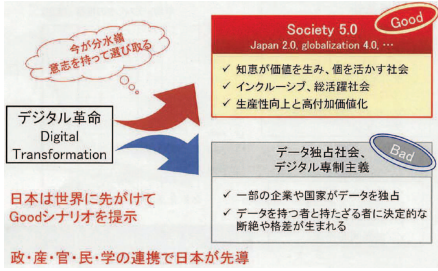
- RISC-V
- RISC-V FoundationというNPOが立ち上げられている



データ・ドリブン

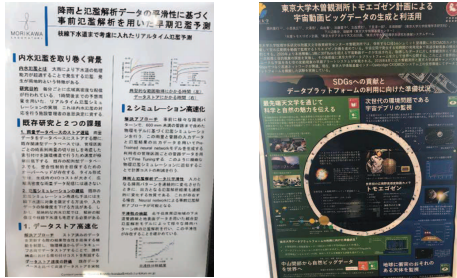
- データとは?
 - データソース
 - 組織や社会に蓄積されたデータ (ビッグデータ)
 - 社会環境や自然環境のセンシングデータ (IoTデータ)
 - データ形式
 - 整理されたデータ (データベース)
 - 生データ
 - データ参照方法
 - データベース検索
 - 検索エンジン
 - AI処理で情報抽出

● データが価値の原点



- ダウンサイジング
 - IBM PC (1981), Apple iPhone (2007)
- Moore's law
 - 「半導体の集積率は18か月で2倍になる」(1965)
- フォンノイマン・コンピュータ
 - EDVACに関する報告書の第一草稿 (1945)
- そろそろ変革が必要？
 - 消費電力と高集積コスト
 - 古くなったから、という訳ではない
 - T型フォードは1908年から
 - プログラミングモデルの議論が不可欠
 - アプリの議論が不可欠

● データ活用の提案例 - 東大で開催された「データ活用社会創成シンポ」より



ありがとうございました

- 現状のAI処理
 - パターン・画像認識の機械学習化に成果
 - 教師有り学習だけでなく、教師無し学習にも成果
 - 膨大なトレーニングデータが必要
 - 結果の説明性について課題 (人間の認識とのマッチング)
 - 仮説の設定、原理の抽出は課題

8. まとめ