

AWS オンラインセミナー

# 夏の Amazon EC2 祭り 2023 最新 インスタンス活用編

～新しいインスタンスを使ってコスト最適化～

**AMD** 

**together we advance\_data centers**

# 注意事項

本プレゼンテーションには、Advanced Micro Devices, Inc. (AMD) の将来の予測に関する記述が含まれています。これには、AMD 製品の特徴、機能性、パフォーマンス、入手可能性、時期、および期待される利益に加えて、AMD 製品のロードマップなどがあり、これらは米国の1995年私的証券訴訟改革法 (Private Securities Litigation Reform Act of 1995) のセーフハーバー条項に従って作成されています。将来の予測に関する記述には、一般的に、「するであろう」、「する可能性がある」、「期待される」、「考えられる」、「計画している」、「意図している」、「予測される」といった用語が含まれます。投資家の皆様は、本プレゼンテーションの将来の予測に関する記述が、現在の判断と仮定・予想に基づき、本プレゼンテーション実施日におけるものであり、実際の結果が現在の予想と大きく異なる可能性のあるリスクと不確実性が含まれていることにご注意ください。当該記述は、既知および未知のリスクや不確実性の影響を受けるものであり、その多くは予測が困難で、一般的に AMD がコントロール可能な範囲を超えるものであるため、実際の結果やその他の将来の事象は、将来の予測に関する情報や文章で表現または暗示・予測されている内容とは大きく異なる可能性があります。投資家の皆様には、Form 10-K および 10-Q に関する AMD の最新の報告書を含め、AMD が証券取引委員会に提出する年次報告書に含まれるリスクと不確実性を詳細に検討することが求められますが、検討すべき事項はこれらに限定されません。

AMD は法律で義務付けられている場合を除き、本プレゼンテーションで発表された将来の予測に関する記述を更新する義務を負わないものとし、本ドキュメントにより当社は当該責任を否認します。

# AGENDA

テクノロジートレンド

ロードマップ

AMD EPYC™ 9004 プロセッサ  
Zen 4 マイクロアーキテクチャー概要

まとめ



# AMD EPYCとAMD INSTINCT MI200アクセラレーター搭載の FRONTIERスーパーコンピューターが TOP500で世界1位をキープ (AS OF MAY 2023)



 **OAK RIDGE** | LEADERSHIP  
National Laboratory | COMPUTING  
FACILITY

**AMD**  
**EPYC**

**INSTINCT**

スーパーコンピューター「トップ500」のうち、  
最もシェアを拡大したCPUとGPUがAMD製品

AMD EPYC™ は日常で使われています  
THE WORLD RELIES ON  
AMD EPYC™ CPUs EVERY DAY



zoom

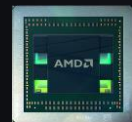
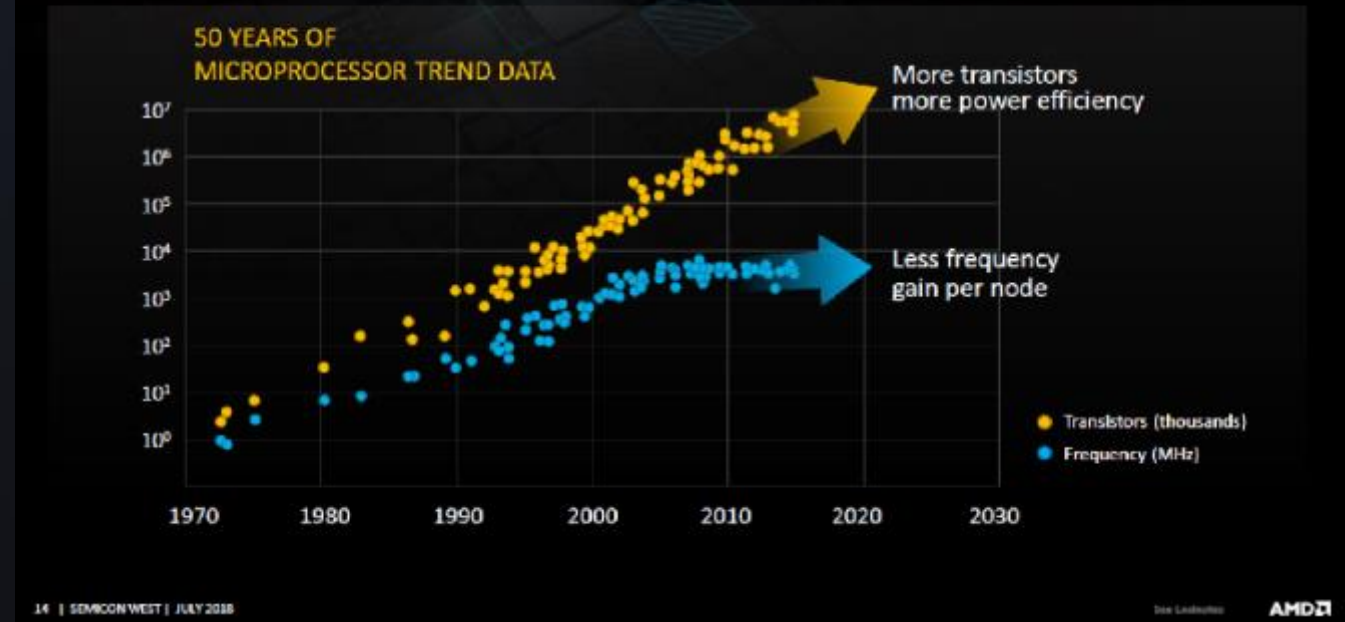
---

マイクロソフト オフィス365, TEAMS, テンセント ミーティング,  
ツイッター, ZOOM そのほか

# 最新 CPU トレンド

- アメリカ CHIPS法案成立
- 世界的な半導体供給不足解消へ
- ムーアの法則が減速
- モノシリック・ダイからチップレットへ
- Arm の台頭
  - Apple, AWS Graviton, 富岳
- 3D パッケージ技術の登場
- 競争が激化しています

## ムーアの法則が減速



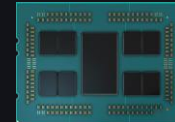
2.5D HBM

2015



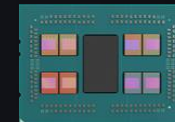
Highest Capacity  
FPGA 2.5D

2017



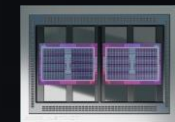
Heterogeneous  
Chiplets

2019



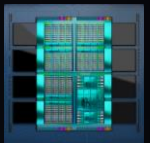
AMD 3D V-Cache™

2021



MI200  
2.5D EFB

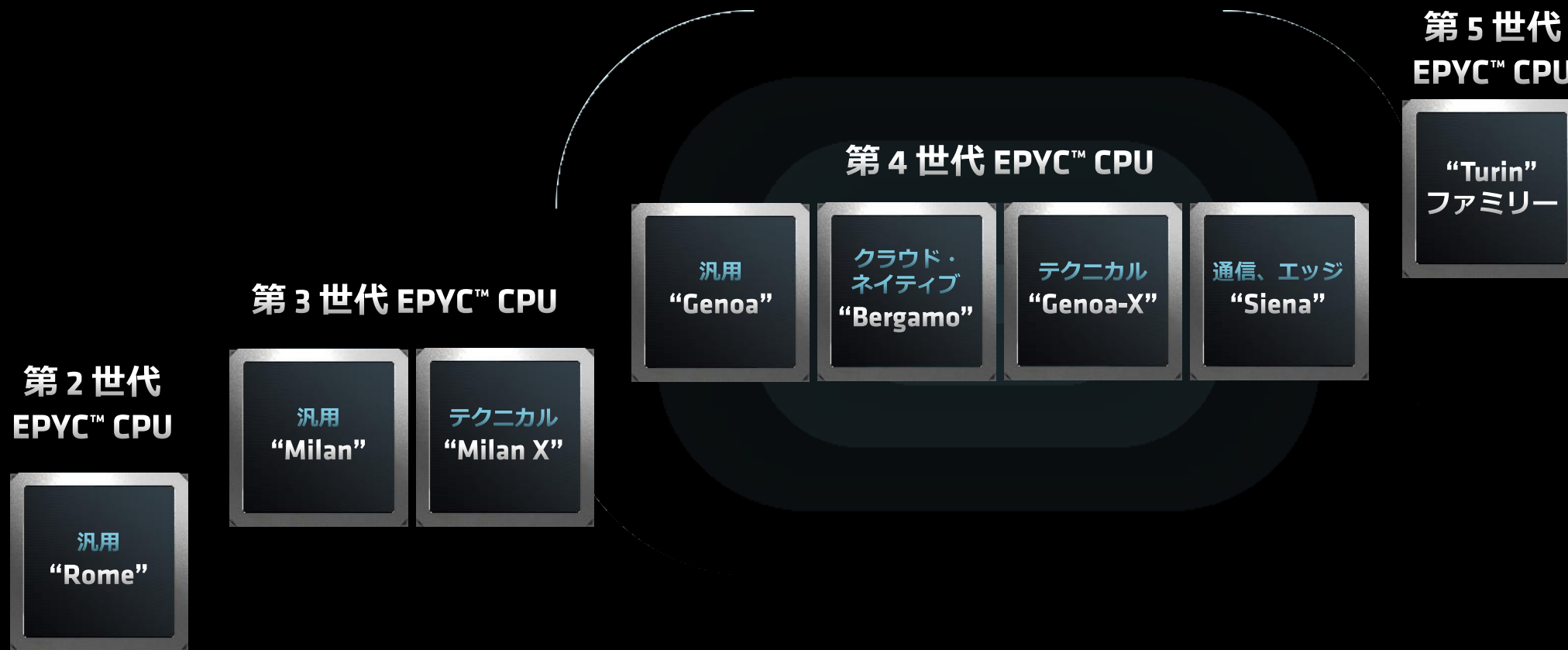
2022



MI300  
3D/2.5D Hybrid

2023

# 戦略とロードマップを確実に実行



2019

2024

ロードマップは予告なしに変更される可能性があります。

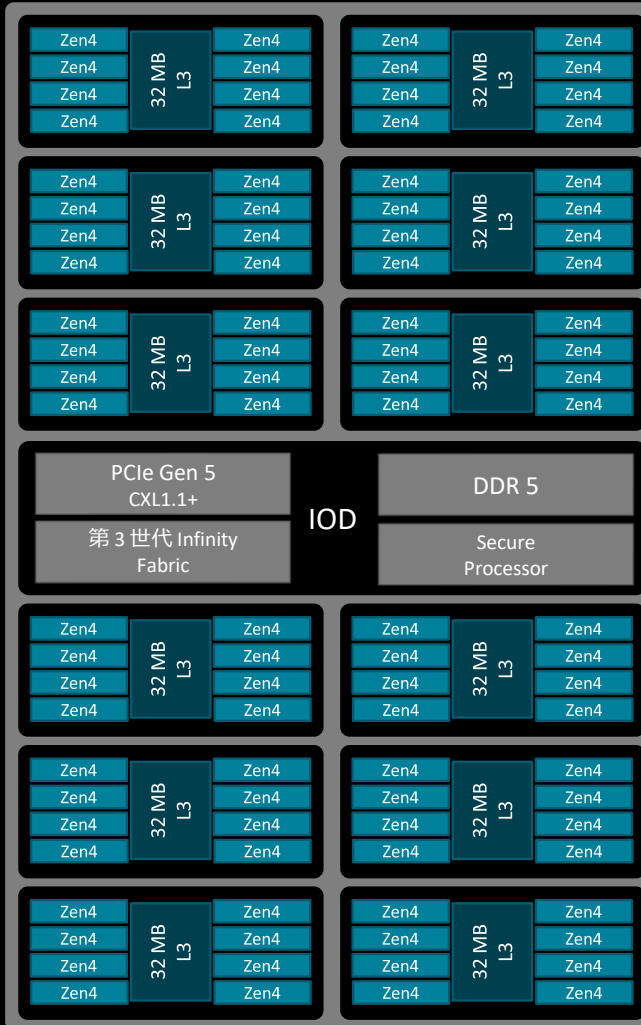
# AMD EPYC™ 9004 概要 (Genoa)

## 演算

- **AMD “Zen4” x86 コア**  
(最大 12 CCD / 96 コア / 192 スレッド)
- 1MB L2/コア、最大 32MB L3/CCD
- **命令セット : bfloat16、VNNI、AVX-512 (256bit データパス)**
- メモリーは 57b/52b 仮想 / 物理アドレスに対応可能
- 最新の IOD と内部の AMD Gen3 Infinity Fabric™ アーキテクチャーに加えて、ダイ・ツー・ダイの帯域幅を拡大
- ターゲット TDP 範囲 : 最大 400W (cTDP)
- **最新の RAS機能**

## メモリー

- **12 チャンネル DDR5 と最大 4800 MHz の ECC**
- オプションで 2、4、6、8、10、12 チャンネルメモリーインターリーブ<sup>1</sup>
- RDIMM、3DS RDIMM
- 最大 2 DIMM/チャンネル容量、2 ソケットシステム (256GB 3DS RDIMM) で最大 12TB<sup>1</sup>



## SP5 プラットフォーム

- **新ソケット、パワーと VR を向上**
- 最大 4 リンクの第 3 世代 AMD Infinity Fabric™ と最大 32Gbps のスピード
- 柔軟なトポロジーオプション
- サーバーコントローラーハブ 内蔵(USB、UART、SPI、I2C など)

## I/O 統合 - チップセットなし

最大 160 IO レーン (2P) の PCIe® Gen5

- 最大 32 Gbps のスピード、x1 までの分岐に対応
- 2P 構成で最大 12 の追加 PCIe® Gen3 レーン (8 レーン - 1P)
- SATA で最大 32 の IO レーン
- CXL1.1+ では 64 I/O レーンをサポート、x4 までの分岐に対応

## セキュリティー機能

専用セキュリティー・サブシステムの強化

Secure Boot、ハードウェアの信頼の基点 (Root of Trust)

SME (Secure Memory Encryption)

SEV-ES (Secure Encrypted Virtualization & レジスタ暗号化)

**SEV-SNP (Secure Nested Paging)、AES-256-XTS と VM 暗号化を向上**

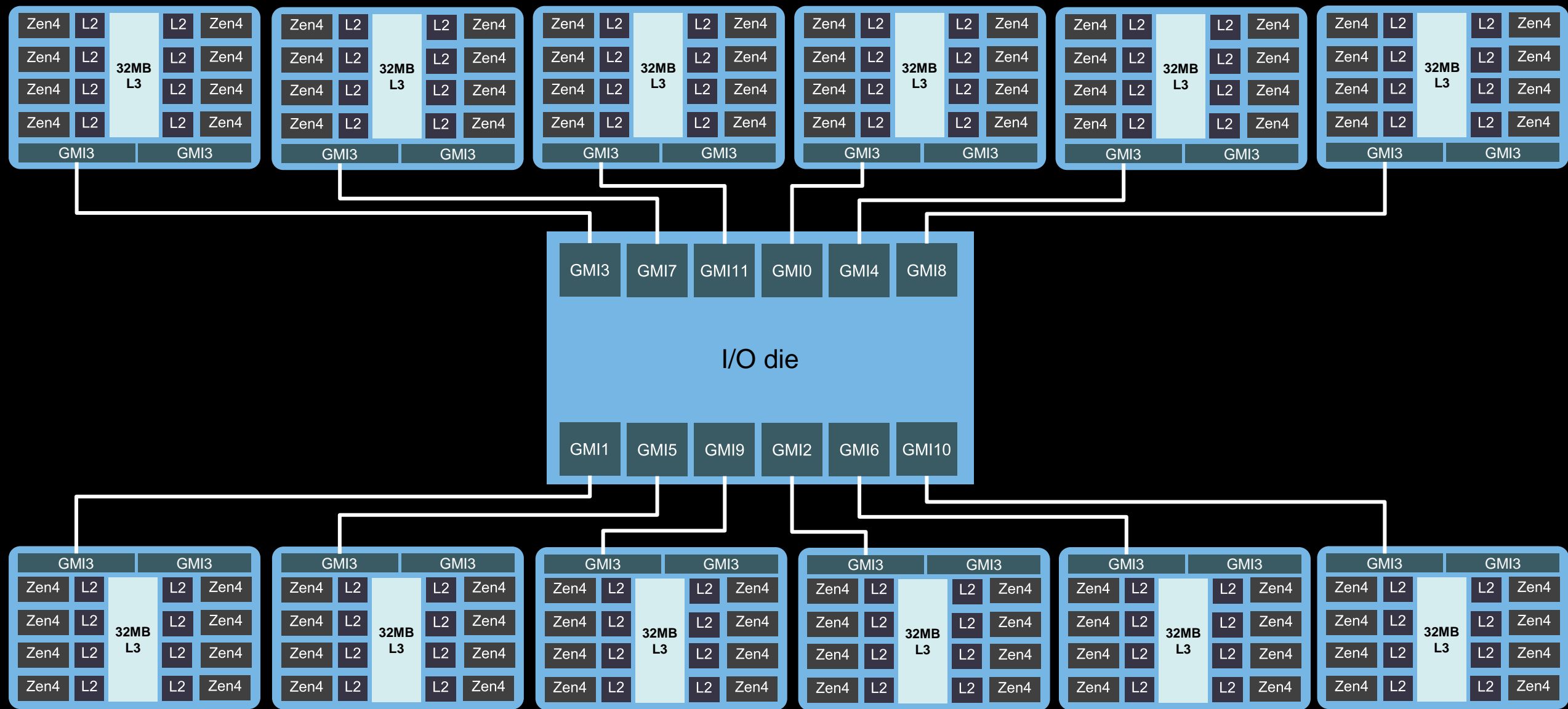
青い文字は前世代CPU、EPYC 7003 からのアップグレードを示します。

<sup>1</sup> 所定の DIMM ボビュレーションルールを適用。GD-183: AMD Infinity Guard の機能は、EPYC™ プロセッサの世代によって異なります。

Infinity Guard のセキュリティー機能は、サーバー OEM および/またはクラウド・サービス・プロバイダーにより有効化されている必要があります。Infinity Guard の詳細については、<https://www.amd.com/en/technologies/infinity-guard> をご覧ください。

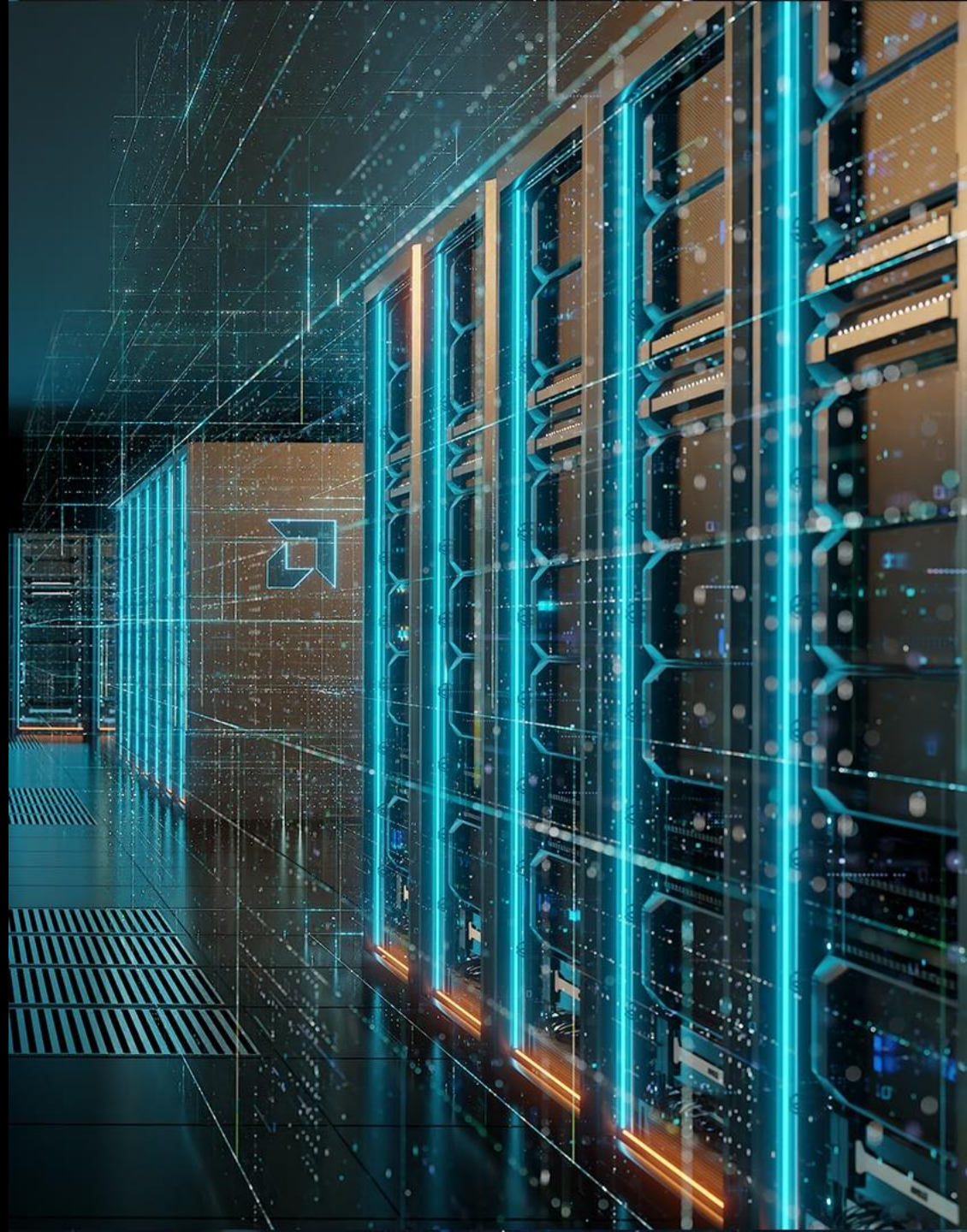


# EPYC 9004 – 12 CCD コンフィギュレーション





# “Zen 4” アーキテクチャー概要



# “Zen 4” アーキテクチャー の設計目標

## パフォーマンス

IPC と周波数の 2 桁向上

## レイテンシー

大容量 L2 による平均レイテンシー低下、  
キャッシュの効率向上

## スループット

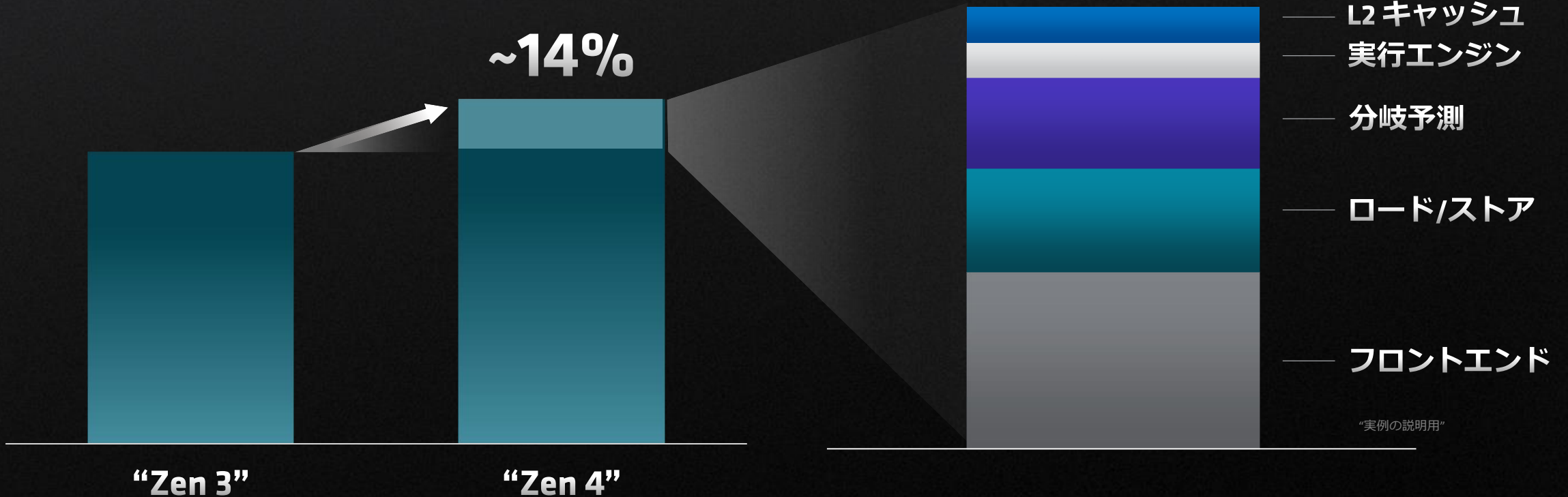
動的電力消費の大幅な削減、コア数の  
効率的な増加

# 世代間の性能向上

“Zen 4” はサーバー CPU の IPC パフォーマンスを ~14% 向上<sup>1</sup>

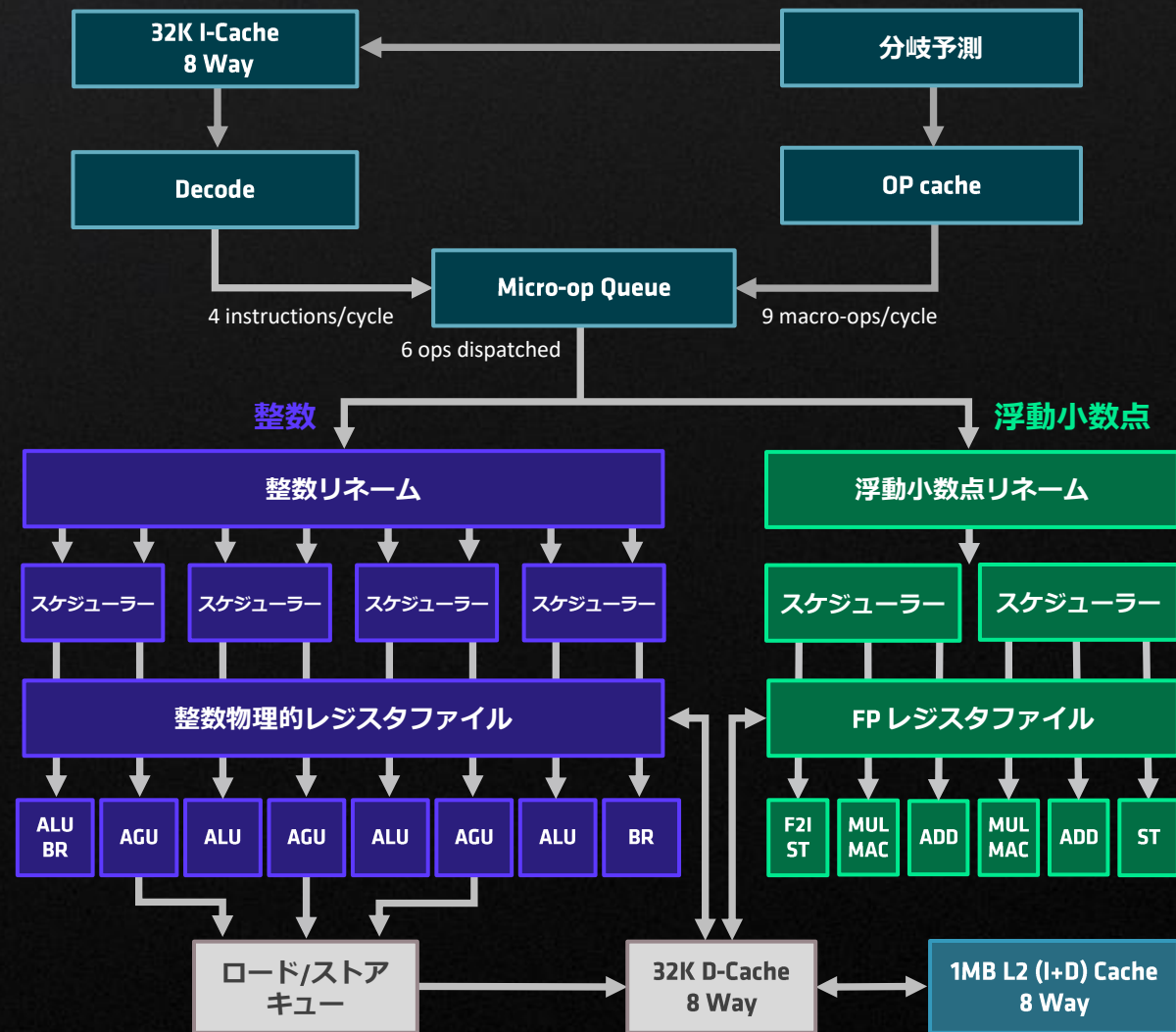
33 種類のサーバーワークロードの幾何平均  
(固定周波数、8+1 CCD)

“Zen 4” パフォーマンスを発揮する理由



# “Zen 4” マイクロアーキテクチャー概要

- “Zen 3” マイクロアーキテクチャーに基づく
- 分岐予測の向上
  - サイクルあたり2つの分岐を予測
  - 大容量 L1 BTB
  - 大容量 L2 BTB
- OP cache 大型化
- 命令リタイヤキュー大型化
- Int/FP レジスタファイル大型化
- コア全体でよりディープなバッファ
- 電力効率に優れた AVX-512 を浮動小数点演算ユニットでサポート
  - 256bit データパス
- ロード/ストアの改善
  - ポートの競合を低減
  - 大容量 L2 DTLB
- L2 Cache 1M、8 Way
- 大容量共有 L3 cache



# “Zen 3” から “Zen 4” への進歩

	“Zen 3”	“Zen 4”
LDQ	72	<b>88</b>
STQ	64	64
Mirco-op cache	4k OP	<b>6.75k OP</b>
L1 I/D-cache	32/32k	32/32k
L2 cache	512k	<b>1M</b>
L3 cache/core	4M	4M
L2 TLB	2k	<b>3k</b>
L2 latency	12 cycles	14 cycles
L3 latency	46 cycles	50 cycles

	“Zen 3”	“Zen 4”
Issue width (Int + FP/SIMD)	10+6	10+6
Int reg	192	<b>224</b>
Int scheduler	96	96
FP reg	160	<b>192</b>
ROB	256	<b>320</b>
FADD/FMUL/FMA latency	3/3/4 cycles	3/3/4 cycles
L1 BTB*	2 x 1k	<b>2 x 1.5k</b>
L2 BTB*	2 x 6.5k	<b>2 x 7k</b>

# AVX-512

## 拡張機能に対応

- AVX512F - ファウンデーション (基盤)
- AVX512DQ - パック整数計算命令
- AVX512\_IFMA - 整数融合積和
- AVX512CD - ループのベクトル化で競合を検出
- AVX512BW - さらにパック整数命令を追加
- AVX512VL - 新しい命令を 128 & 256 bit に拡張
- AVX512\_VBMI - ベクトルバイトの順列処理
- AVX512\_VBMI2 - ベクトルバイトの順列処理を向上
- GFNI - ガロア体の新しい命令 (SSE、AVX、AVX512)
- AVX512\_VNNI - ベクトル NN 命令
- AVX512\_BITALG
- AVX512\_VPOPCNTDQ
- AVX512\_BF16 - BFloat16 変換

## Neural Magic Sparse INT8 推論

~4.2倍

NLP スループット

2P AMD EPYC™ 9654 と 2P AMD EPYC™ 7763 の比較

~3倍

画像分類スループット

2P AMD EPYC™ 9654 と 2P AMD EPYC™ 7763 の比較

~3.5倍

オブジェクト検出スループット

2P AMD EPYC™ 9654 と 2P AMD EPYC™ 7763 の比較

# AVX-512 インストラクションセットのサポート詳細

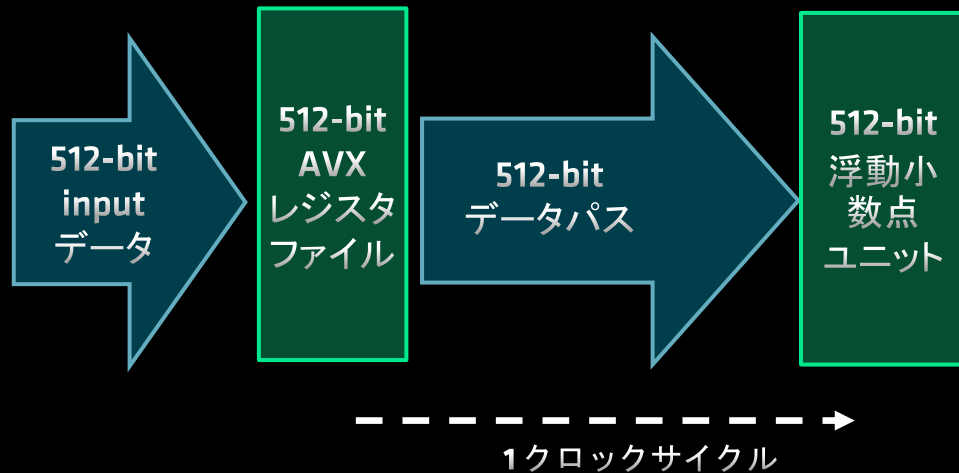
AVX-512 Subset	F	CD	ER	PF	4FMAPS	4VNNIW	VPOPCNTDQ	VL	DQ	BW	IFMA	VBMI	VNNI	BF16	VBMI2	BITALG	VPCLMULQDQ	GFNI	VAES	VP2INTERSECT	FP16	
“Skylake”	Y								Y								N					
“Cascade Lake”	Y						N		Y			N	Y					N				
“Cooper Lake”	Y								Y			N	Y	Y				N				
“Ice Lake”	Y				N				Y					N			Y				N	
“Genoa”	Y				N				Y								Y				N	
“Sapphire Rapids”	Y				N				Y								Y				N	Y

AMD “Zen 4” AVX-512 インストラクションは第3世代 Intel Xeon (“Ice Lake”)とFP16を除いた第4世代 Intel Xeon (“Sapphire Rapids”)と命令セットの互換性があります。

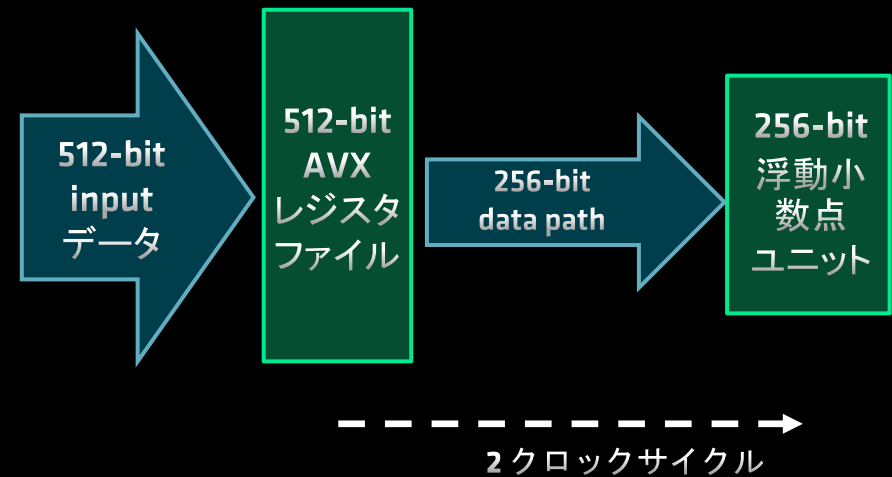


# Intel AVX-512 vs AMD “Zen 4” AVX-512

- Intel データパス幅 512 ビット

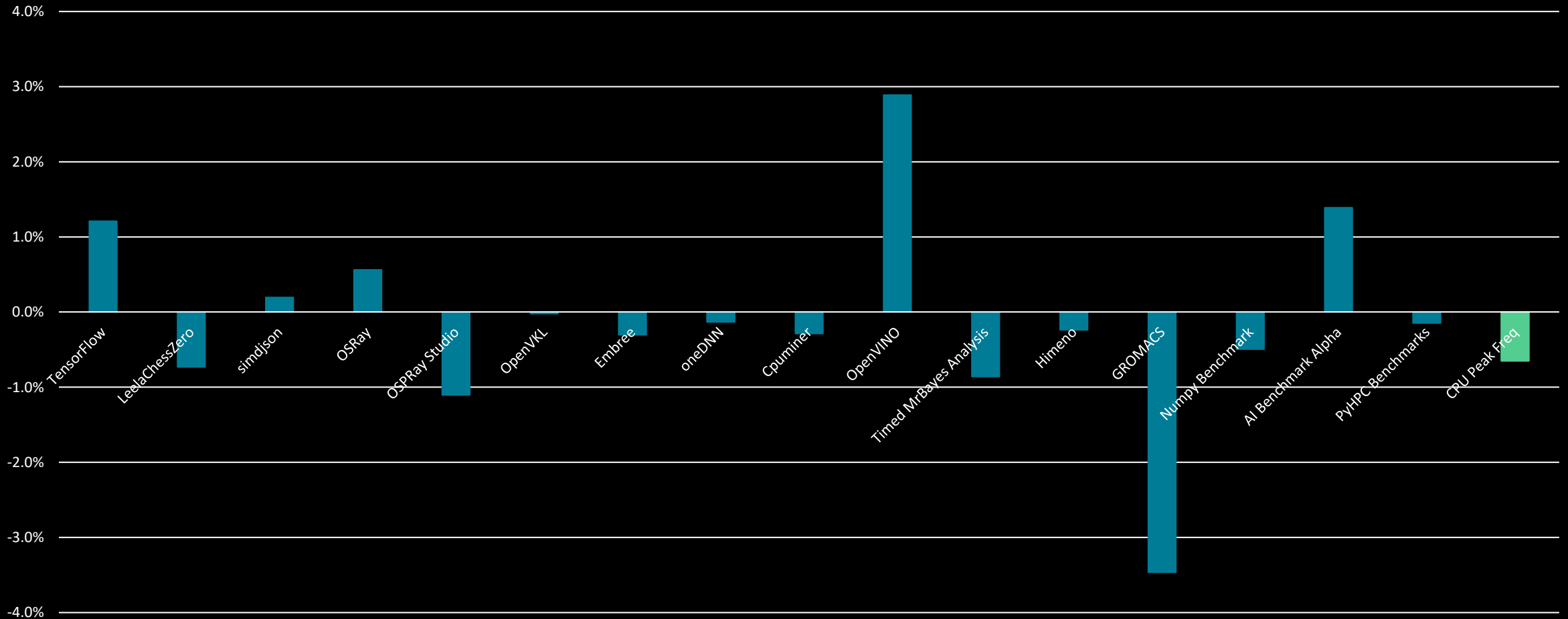


- “Zen 4” データパス幅 256 ビット



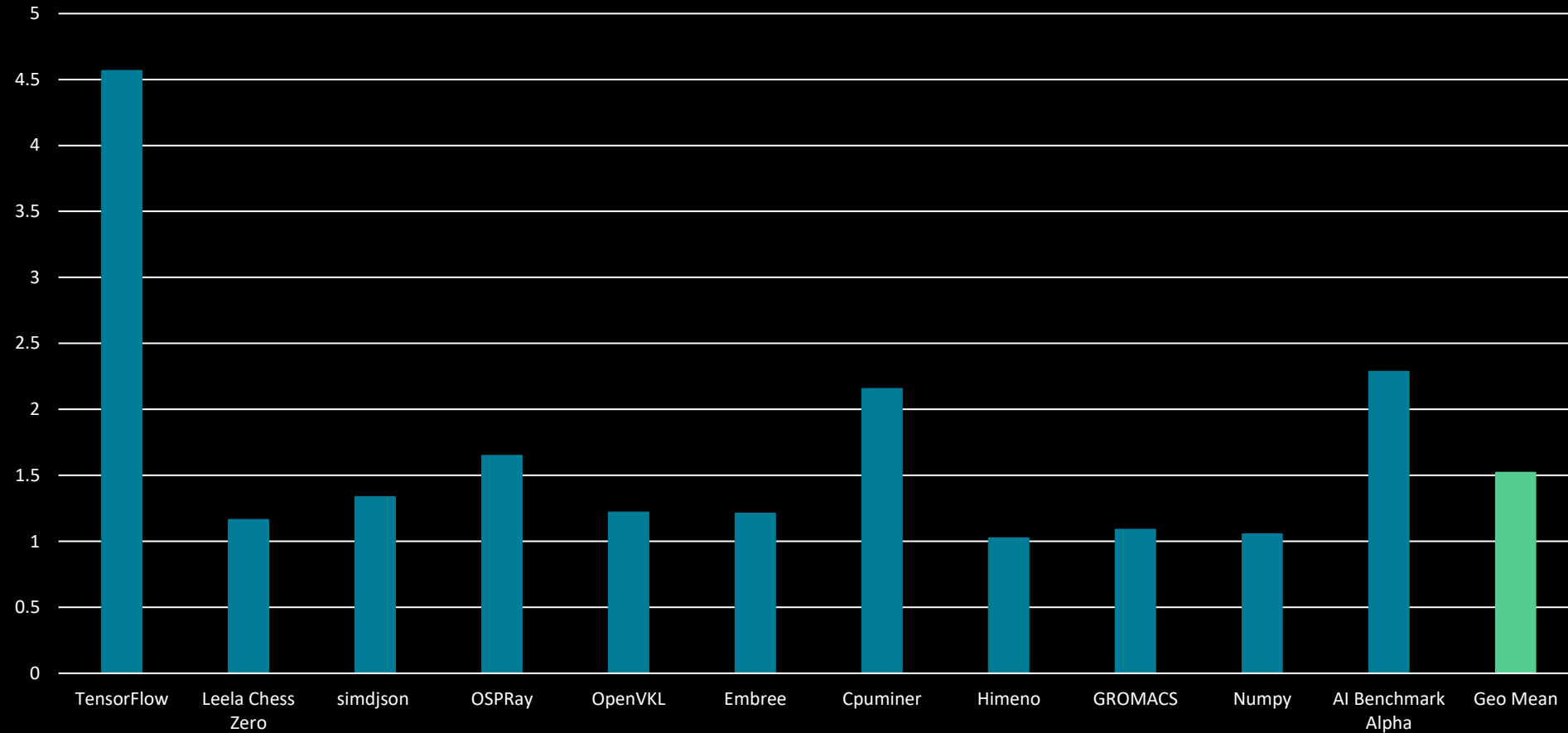
- AMD AVX-512は小面積で高効率な実装を目指しました
- AMD AVX-512 デザインは、CPUコアと同じ256ビットのデータ・パスを使用し、2つのパーツがシーケンシャルなクロック・サイクルで実行できるように設計している
- “Zen 4” コアは、競合のGolden Coveコアと比較して約50% 小面積で約 47% 高効率 EPYC-041

# EPYC™ 9654: AVX-256と比較したAVX-512の相対コア周波数 - Phoronix\*\*



- AI/MLで利用される命令（AVX-512 VNNI + フロントエンドの最適化）では効率的な実行が可能で、特定のワークロードで周波数が向上
- AVX-512の平均周波数はAVX-256とほぼ同じ(-0.7%)

# EPYC™ 9654: AVX-256と比較したAVX-512の相対性能/ワット - Phoronix\*\*



- 効率的なAMD AVX-512の実装は、AVX-256と比較して、TensorFlowなどのAI/MLワークロードの性能/ワットを大幅に向上させます。

# GenoaのAVX512について

## Zen4における

- AVX512の実装は、AVX256と比較して35%性能が向上しました
- AVX512の平均動作周波数はAVX256とほぼ同じ (-0.7%)
- AVX512の実装はエネルギー効率に優れており、TensorFlowなどのAI/MLワークロードで性能/ワットで大きな利点があります。

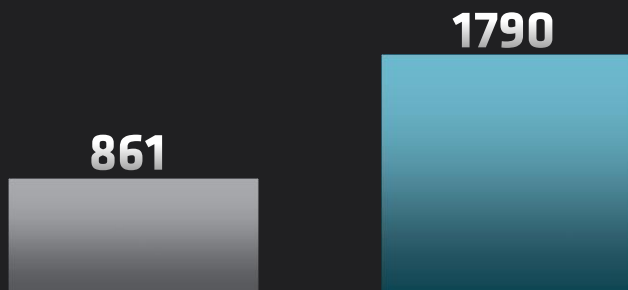


# 業界をリードする先進的な AMD EPYC™ CPU

## クラウド

SPECrate®2017\_int\_base

**107%** より高いパフォーマンス



第3世代 EPYC™  
**7763**

第4世代 EPYC™  
**9654**

整数演算スループット  
クラウド・サービス・プロバイダー

IaaS / PaaS

検索

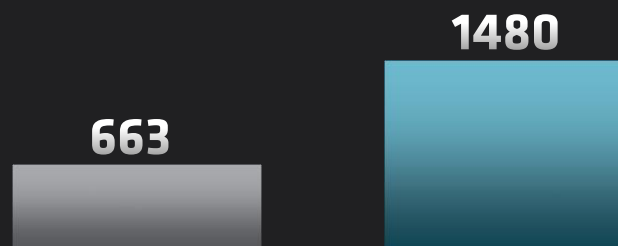
ソーシャル

SaaS

## HPC

SPECrate®2017\_fp\_base

**123%** より高いパフォーマンス



第3世代 EPYC™  
**7763**

第4世代 EPYC™  
**9654**

浮動小数点演算スループット  
ハイパフォーマンス・コンピューティング

設計 &  
シミュレーション

研究 &  
学術研究

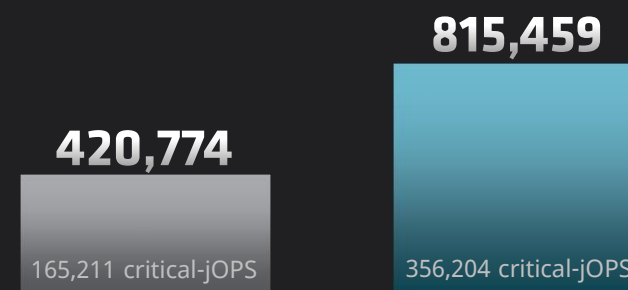
機械学習

スーパー  
コンピューティング

## エンタープライズ

SPECjbb®2015 MultiJVM max-jOPS

**94%** より高いパフォーマンス



第3世代 EPYC™  
**7763**

第4世代 EPYC™  
**9654**

サーバーサイド Java® 最大スループット  
エンタープライズ IT

仮想化

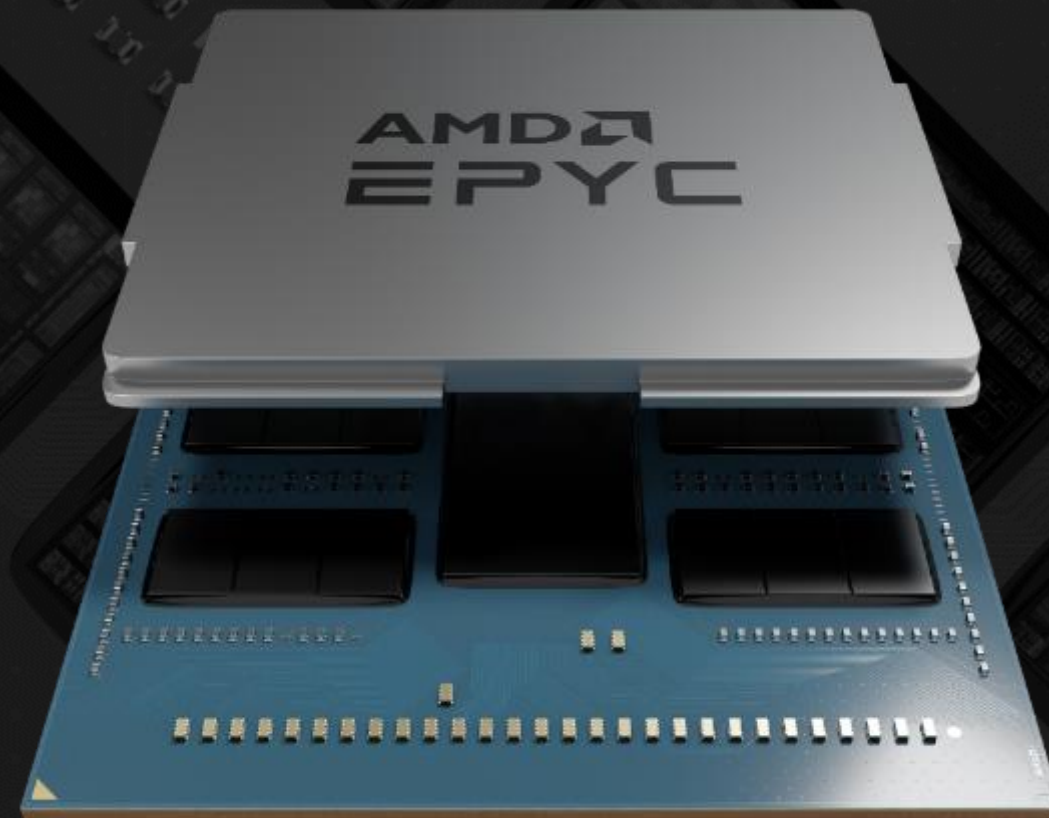
SDS/HCI

Hadoop

NoSQL

# AMD EPYC™ 9004 シリーズ CPU

- **Zen4 でコア・パフォーマンスを大幅にアップグレード**
  - ~14% IPC、AVX-512 ISA サポート (VNNI/bfloat16、1 MB/c L2\$ などを含む)
- **CPU あたり最大 96 コア**
- **ハイブリッド型チップレット・アーキテクチャー、最先端の 5nm + 6nm プロセス**
- **リーダーシップのプラットフォーム機能、スケーリングに対応したプロビジョニング**
  - シングルランクのパフォーマンスを強化した 12 チャンネル DDR5-4800
  - 128L 32Gbps + 8L 8Gbps 多機能 SERDES
- **CXL1.1+ メモリーサポート (CXL™ “タイプ 3”)**
  - 高度なメモリーアタッチ機能で DDR および新規規格のメモリーに対応
  - SEV-SNP、QoS、階層化メモリー管理拡張
- **サーバー RAS : データポイズニング、AMD-C ECC、帯域外の重大なエラーを収集**
- **SEV-SNP キー拡張によるセキュリティ強化、AES-256-XTS サポート**



**AMD** 



# Useful link

AMD EPYC ホワイトペーパー

<https://www.amd.com/en/processors/server-tech-docs/search>

AMD EPYC パフォーマンス・チューニングガイド

<https://www.amd.com/en/processors/tuning-guides-server>

AMD EPYC デモ

<https://www.amd.com/ja/processors/epyc-9004-demo>

AMD EPYC ツール (TCO、CPUセレクトター)

<https://www.amd.com/ja/processors/epyc-tools>

# ENDNOTES

EPYC-38: Based on AMD internal testing as of 09/19/2022, geomean performance improvement at the same fixed-frequency on a 4th Gen AMD EPYC™ 9554 CPU compared to a 3rd Gen AMD EPYC™ 7763 CPU using a select set of workloads (33) including est. SPECrate®2017\_int\_base, est. SPECrate®2017\_fp\_base, and representative server workloads.

EPYC-041: ~40% less area measures Core + L2 Area: “Zen 4” = 3.84 mm<sup>2</sup> vs. “Sunny Cove” ~6.5mm<sup>2</sup>. ~48% estimated SPECrate®2017\_int\_base GCC -O3 Jemalloc configured estimate based on internal AMD reference platform and Intel platform measurements of 11/10/2022. Comparison of estimated 1P AMD EPYC 9534 (537 est. SPECrate®2017\_int\_base, set to 270 Total TDP W, 64 Total Cores, AMD Est) is 1.48x the performance per watt of 1P Intel Xeon Platinum 8380 (363 est. SPECrate®2017\_int\_base, 270 Total TDP W, 40 Total Cores, AMD est.) for 1.48x the SoC performance/watt. SPEC®, SPEC CPU®, and SPECrate® are registered trademarks of the Standard Performance Evaluation Corporation. See [www.spec.org](http://www.spec.org) for more information. OEM published scores will vary based on system configuration and determinism mode used (claim uses 270W cTDP performance profile)

System configuration details: AMD Quartz CRB BIOS 1002E (NUMA nodes per socket: NPS1, L1 HW Stream Prefetcher: Disabled, L2 HW Stream Prefetcher: Disabled, cTDP=Manual, set to 270w, PPT=Manual, set to 270w, 4-link xGMI speed: 20Gbps, Determinism Control: Manual, set to Performance Determinism), 768GB DDR5-3200 (12x 64GB 2R downclocked from 4800GT/s), Ubuntu 22.04, GCC 10.2.

Intel server platform XYZ BIOS #.# (default BIOS), 512GB DDR4-3200 (16x 64GB), Ubuntu 22.04, GCC 10.2.

DDR	DDR4-3200	DDR5-3200
Core/L3		4MB/core L3
System B/W Triad	321GB/s	345GB/s
CPU / Power	8380 270W	9534 270W
Est SPECrate®2017_int_rate GCC -O3 Jemalloc	363	537
Scale		1.479x

SP5-022: Neural Magic measured results on AMD reference systems as of 9/29/2022. Configurations: 2P EPYC 9654 “Titanite” vs. 2P EPYC 7763 “DaytonaX” running on Ubuntu 22.04 LTS, Python 3.9.13, pip==22.12/deepsparse==1.0.2. BERT-Large Streaming Throughput items/sec (seq=384, batch 1, 48 streams, INT8 + sparse) using SQuAD v1.1 dataset; ResNet50 Batched Throughput items/sec (batch 256, single-stream, INT8 sparse) using ImageNet dataset; YOLOv5s Streaming Throughput ([image 3, 640, 640], batch 1, multi-stream, per-stream latency <=33ms) using COCO dataset. Testing not independently verified by AMD.

PRX-1: “AVX-512 Performance Comparison: AMD Genoa vs. Intel Sapphire Rapids & Ice Lake” - January 2023, page 8, <https://www.phoronix.com/review/intel-sapphirerapids-avx512>

PRX-2: “AMD 4th Gen EPYC 9654 “Genoa” AVX-512 Performance Analysis” - December 2022, page 9, <https://www.phoronix.com/review/amd-epyc-avx512>

# DISCLAIMER AND ATTRIBUTION

The information contained herein is for informational purposes only and is subject to change without notice. While every precaution has been taken in the preparation of this document, it may contain technical inaccuracies, omissions and typographical errors, and AMD is under no obligation to update or otherwise correct this information. Advanced Micro Devices, Inc. makes no representations or warranties with respect to the accuracy or completeness of the contents of this document, and assumes no liability of any kind, including the implied warranties of noninfringement, merchantability or fitness for particular purposes, with respect to the operation or use of AMD hardware, software or other products described herein. No license, including implied or arising by estoppel, to any intellectual property rights is granted by this document. Terms and limitations applicable to the purchase or use of AMD's products are as set forth in a signed agreement between the parties or in AMD's Standard Terms and Conditions of Sale.

Timelines, roadmaps, and/or product release dates shown in these slides are plans only and subject to change. "Rome", "Milan", "Genoa", "Zen", "Zen2", "Zen3", and "Zen4" are codenames for AMD architectures, and are not product names.

©2022 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, EPYC, Infinity Fabric, and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other product names used in this publication are for identification purposes only and may be trademarks of their respective companies. Microsoft SQL Server is a trademark of Microsoft Corporation in the US or other jurisdictions, PCIe - PCIe is a registered trademark of PCI-SIG Corporation, Oracle is a registered mark of Oracle and/or its affiliates, Radioss is a trademark of Altair Engineering, Inc. SPEC®, SPECrate®, and SPECpower® are registered trademarks of Standard Performance Evaluation Corporation. Learn more at [spec.org](http://www.spec.org). More information about SPEC® is available at <http://www.spec.org>. VMmark is a registered trademark of VMware in the US or other countries.