

# Persistent Memory の最新技術動向紹介

今村 智史 † (富士通株式会社)

柿原 俊男 † (HGST Japan)

SNIA-J 次世代不揮発メモリ分科会  
† 会長、‡ 副会長

# 自己紹介（今村 智史）

- 所属：富士通株式会社 ICTシステム研究所
  - メモリ／ストレージシステムに関する研究に従事
  - PMEM の性能分析やユースケース検討
- 専門分野：コンピュータアーキテクチャ
  - OS（特に仮想メモリ周り）も得意です
- 出身：九州大学（博士（工学）取得）
- **SNIA-J 次世代不揮発メモリ分科会 副会長**



# はじめに

- 本講演では Persistent Memory 製品の性能評価結果をご紹介します
- 今回の性能評価実験はすべて柿原さん（HGST Japan）に実施していただき、その実験結果をご提供いただきました

# 目次

- Persistent Memory (PMEM) とは
- 2世代の PMEM 製品
- 性能評価

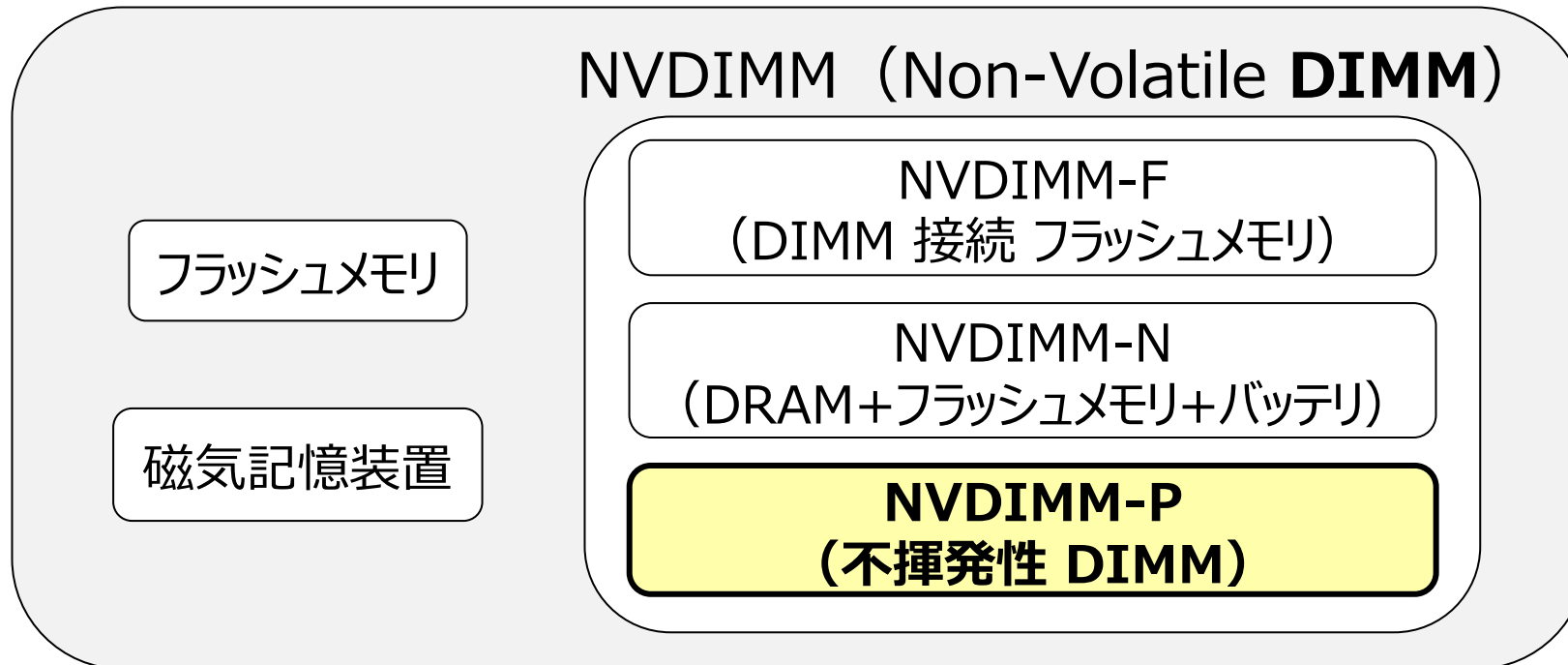
# Persistent Memory (PMEM) とは

# Non-Volatile DIMM (NVDIMM)

“ロード/ストアプログラミングモデルに  
適した性能特性を持つストレージテクノロジー”

出典：[SNIA NVM Programming Model Version 1.2](#)

NVM (Non-Volatile Memory)



# Persistent Memory (PMEM)



出典 : [Intel](#)



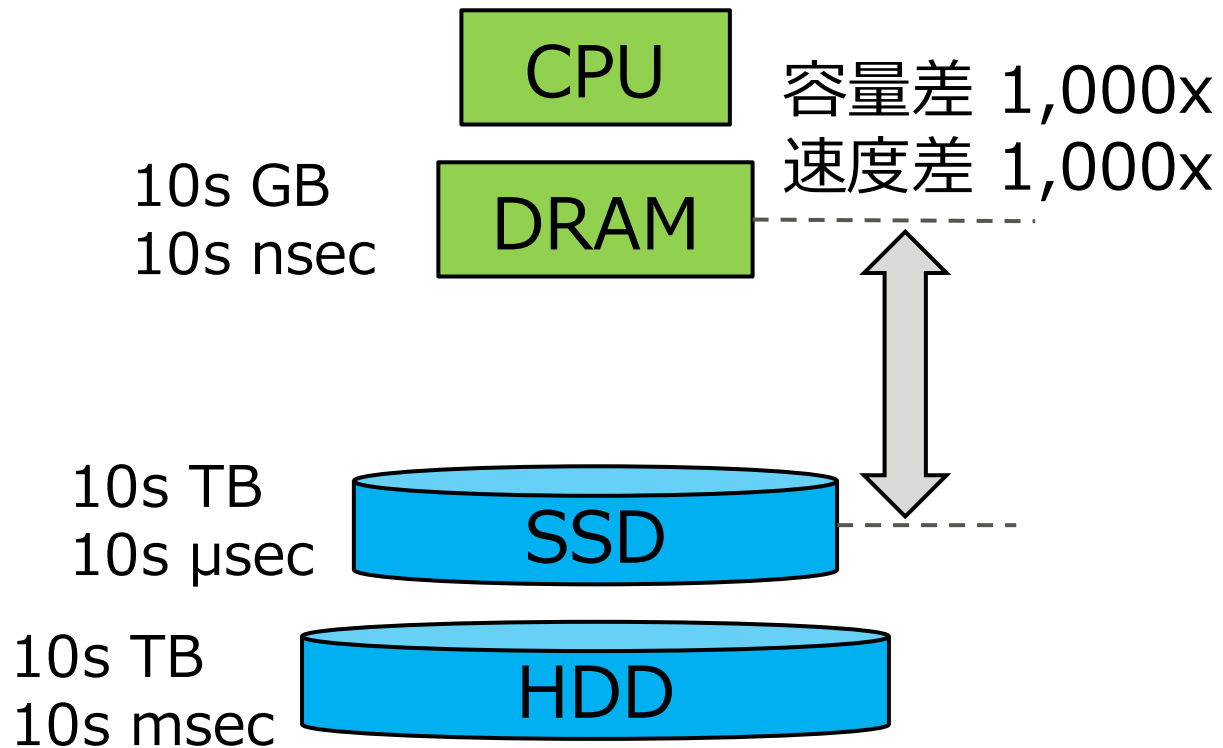
出典 : [Lenovo](#)

- 世界初の NVDIMM-P 製品
  - 3D XPoint™ メモリ技術を採用
- DIMM 当たり最大 **512 GB**
  - CPU あたり最大 4 TB
- DIMM スロットに接続
- PMEM 対応の CPU が必要

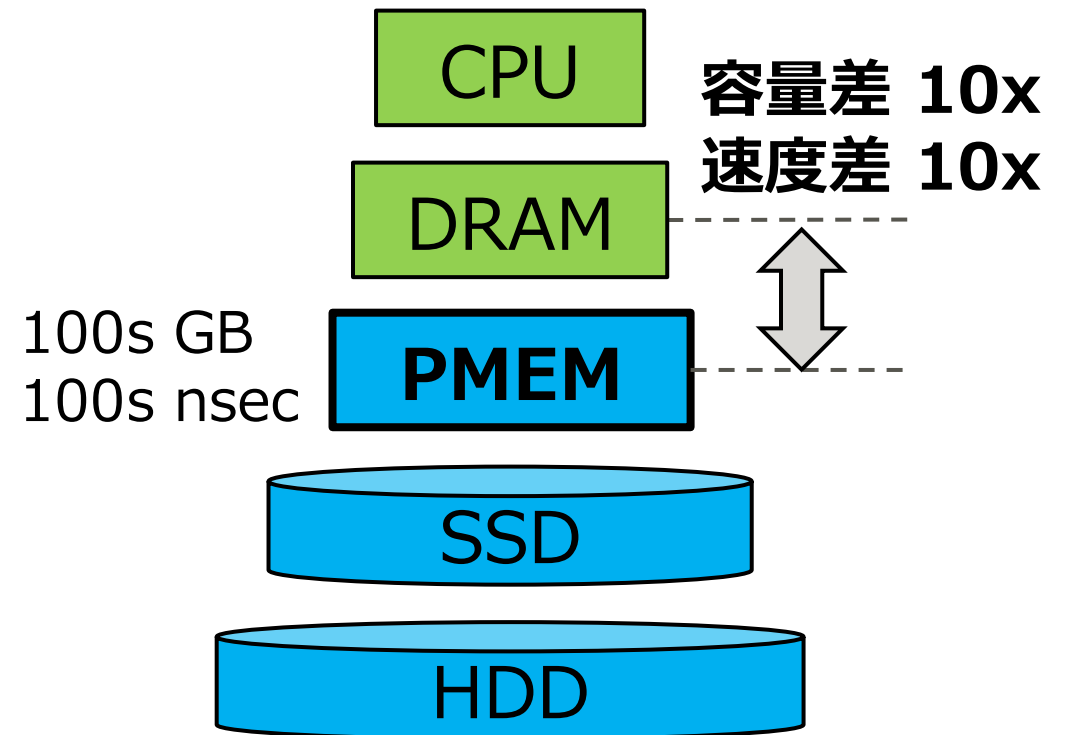
# PMEM の狙い

揮発性  
不揮発性

## 従来のメモリ階層



## 新たなメモリ階層





# PMEM の特徴

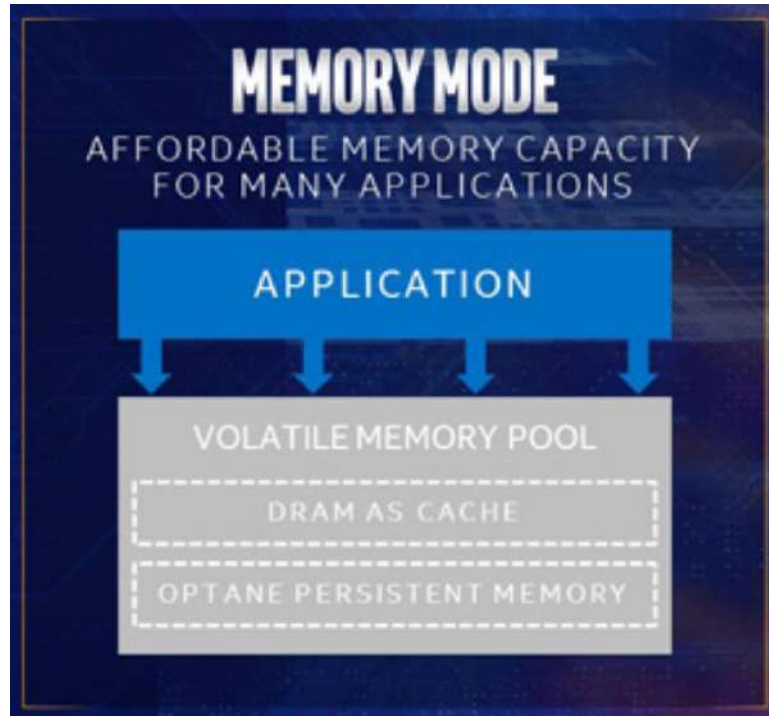
## 利点

- DRAM に比べ容量単価が安く大容量
- バイト単位のロード／ストアアクセスが可能
- 不揮発性を有する（電源遮断時でもデータを保持可能）
- SSD に比べ高性能かつ高書き込み耐性

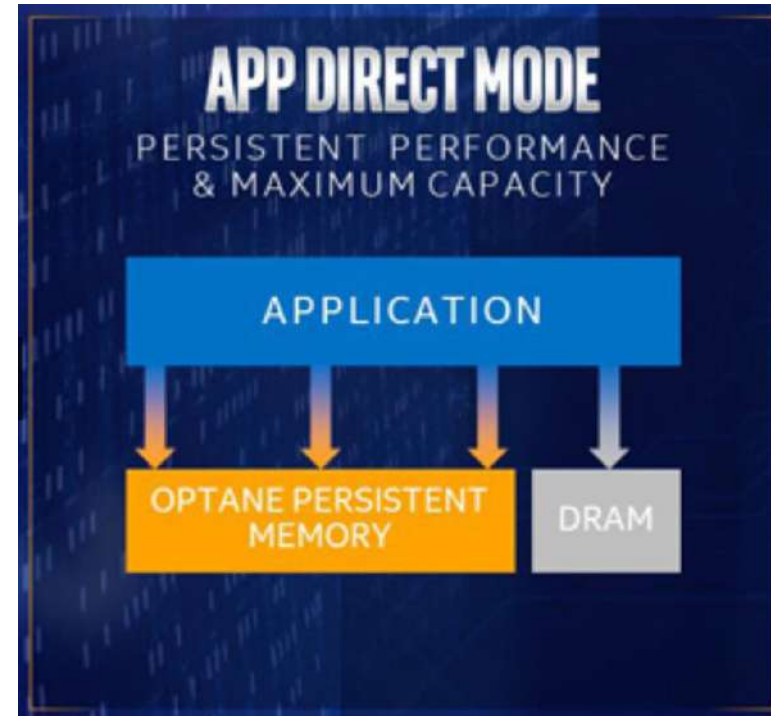
## 欠点

- DRAM に比べ低性能（特に書き込み性能が低い）
- DRAM より低い書き込み耐性

# PMEM の 2 種類のモード



- PMEM を安価で大容量の**揮発性メモリ**として使用
- DRAM はキャッシュとして機能
- アプリの修正必要なし



- DRAM と PMEM を使い分け可能
- PMEM の不揮発性を活用可能
  - PMEM プログラミングが必要

出典：インテル

# PMEM 活用方法の選択肢

選択肢	モード	メリット
DRAM 代替の揮発性メインメモリ	Memory	DRAM と同様に利用可能 (アプリの修正必要無し)
DRAM とは別個の揮発性メインメモリ	App Direct	DRAM/PMEM を容易に使い分け可能
永続ブロックストレージ	App Direct	<ul style="list-style-type: none"><li>• 従来のファイル API が使用可能</li><li>• データを永続化可能</li></ul>
不揮発性メインメモリ	App Direct	メインメモリ上でデータを永続化可能

# 2世代のPMEM製品

# PMEM 100 Series と PMEM 200 Series

## ■ PMEM 100 Series

- 第1世代の Optane™ Persistent Memory (旧 Optane™ DCPM)
- コードネーム : **A**pache Pass
- リリース時期 : Q2'19

## ■ **PMEM 200 Series**

- 第2世代の Optane™ Persistent Memory
- コードネーム : **B**arlow Pass
- リリース時期 : Q2'20

# PMEM 100 Series と 200 Series の比較

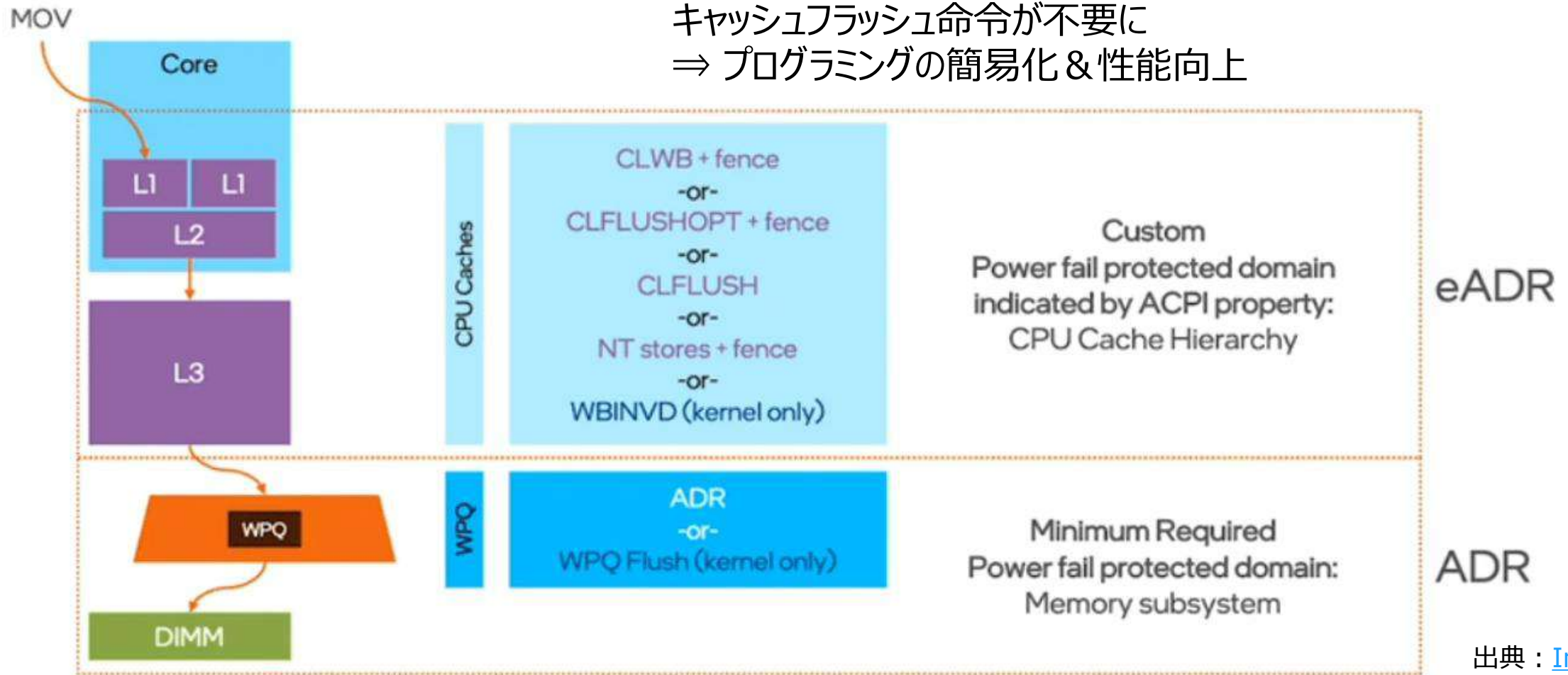
Intel® Optane™ PMem 100 series		Intel Optane PMem 200 series		Intel Optane PMem 200 series	
AES 256-BIT encryption		UP TO 32% higher average memory bandwidth over the previous generation*		AES 256-BIT encryption	
Secure Erase				Secure Erase	
Up To 512 GB modules				Up To 512 GB modules	
Intel® Optane™ PMem 100 series 2nd Generation Intel® Xeon® Scalable processors on 2S/4S/8S platform		Intel Optane PMem 200 series 3rd Generation Intel Xeon Scalable processors on 4S platform		Intel Optane PMem 200 series 3rd Generation Intel Xeon Scalable processors on 2S platform	
8-28 cores	6 channels memory	18-28 cores	6 channels memory	16-40 cores	8 channels memory
3 TB Intel Optane PMem per socket*	4.5 TB Total system memory per socket*	3 TB Intel Optane PMem per socket*	4.5 TB Total system memory per socket*	4 TB Intel Optane PMem per socket**	6 TB Total system memory per socket**
2,666 MT/s DDR4 + Intel Optane PMem		2,666 MT/s DDR4 + Intel Optane PMem		3,200 MT/s DDR4 + Intel Optane PMem	
18 W Max thermal design power		eADR	15 W Max thermal design power	eADR	15 W Max thermal design power

+ Based on testing by Intel as of April 27, 2020 (Baseline) and March 31, 2020 (New).  
 \* 3 TB Intel Optane PMem = 6 x 512 GB Intel Optane PMem per socket, 4.5 TB System Memory = 6 x 512 GB Intel Optane PMem per socket + 6 x 256 GB  
 \*\* 4 TB Intel Optane PMem = 8 x 512 GB Intel Optane PMem per socket, 6 TB System Memory = 8 x 512 GB Intel Optane PMem per socket + 8 x 256 GB

出典 : [Intel](https://www.intel.com)

# eADR (extended Asynchronous DRAM Refresh)

キャッシュフラッシュ命令が不要に  
⇒ プログラミングの簡易化 & 性能向上



出典 : [Intel](#)

# 性能評価

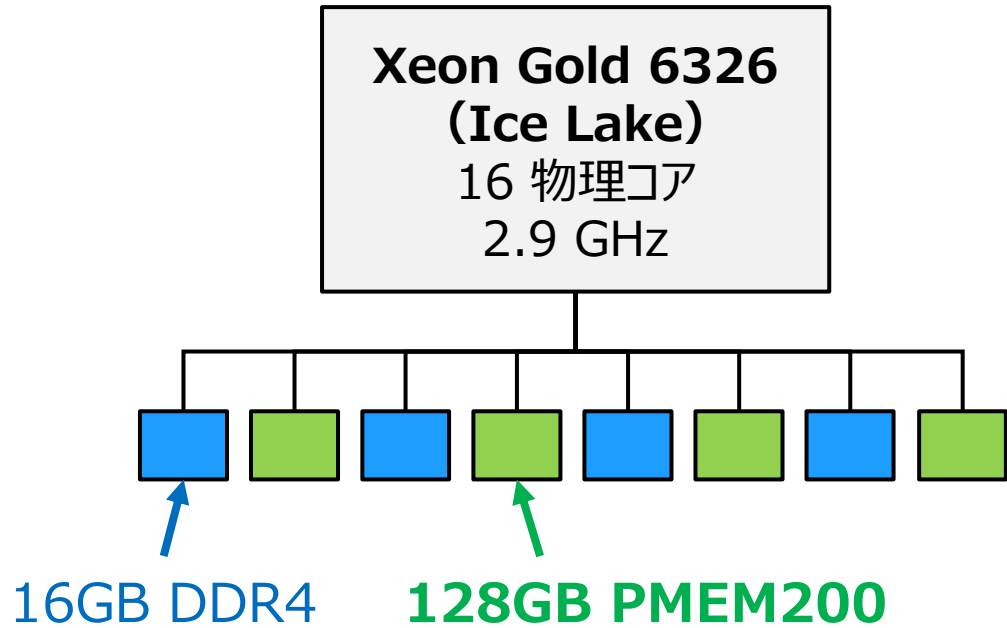


# 評価項目

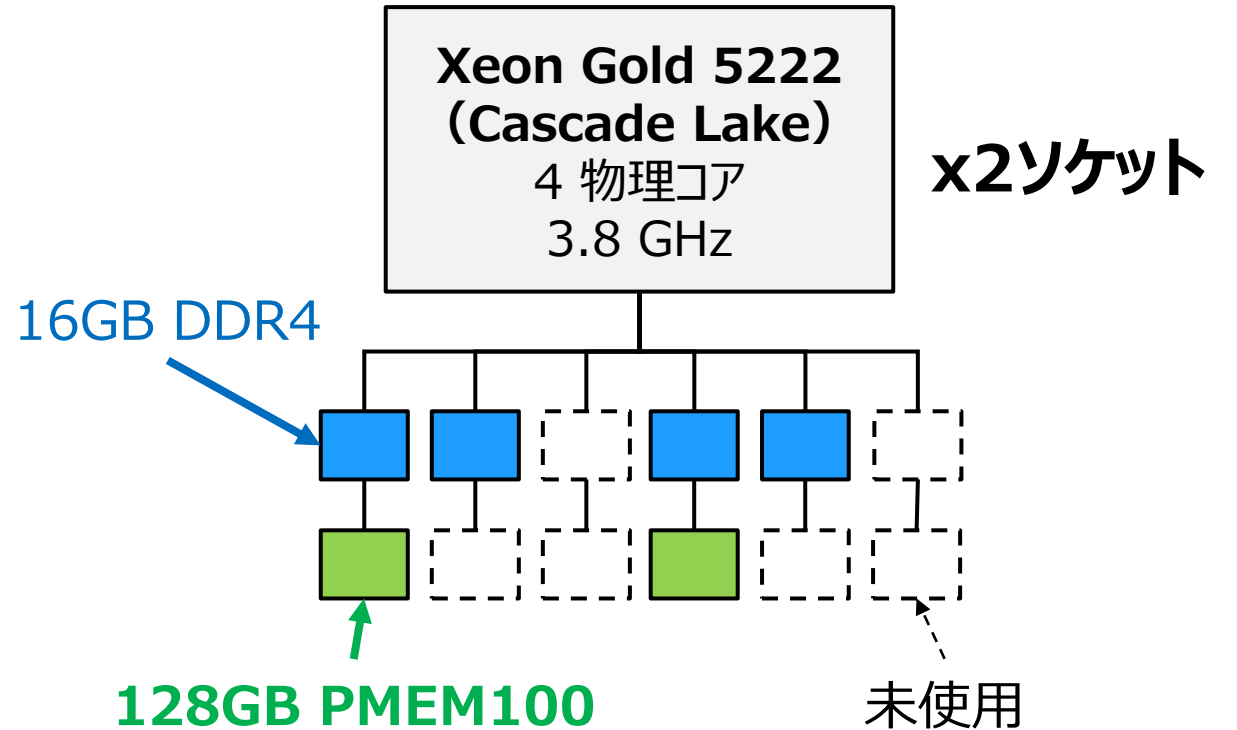
- メインメモリとしての基礎性能
- ブロックデバイスとしての基礎性能
- Spec CPU 2017 ベンチマークの性能
- インメモリキーバリューストアの性能

# 評価環境

## PMEM200 サーバ



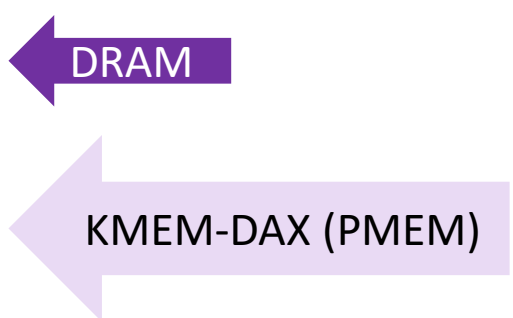
## PMEM100 サーバ



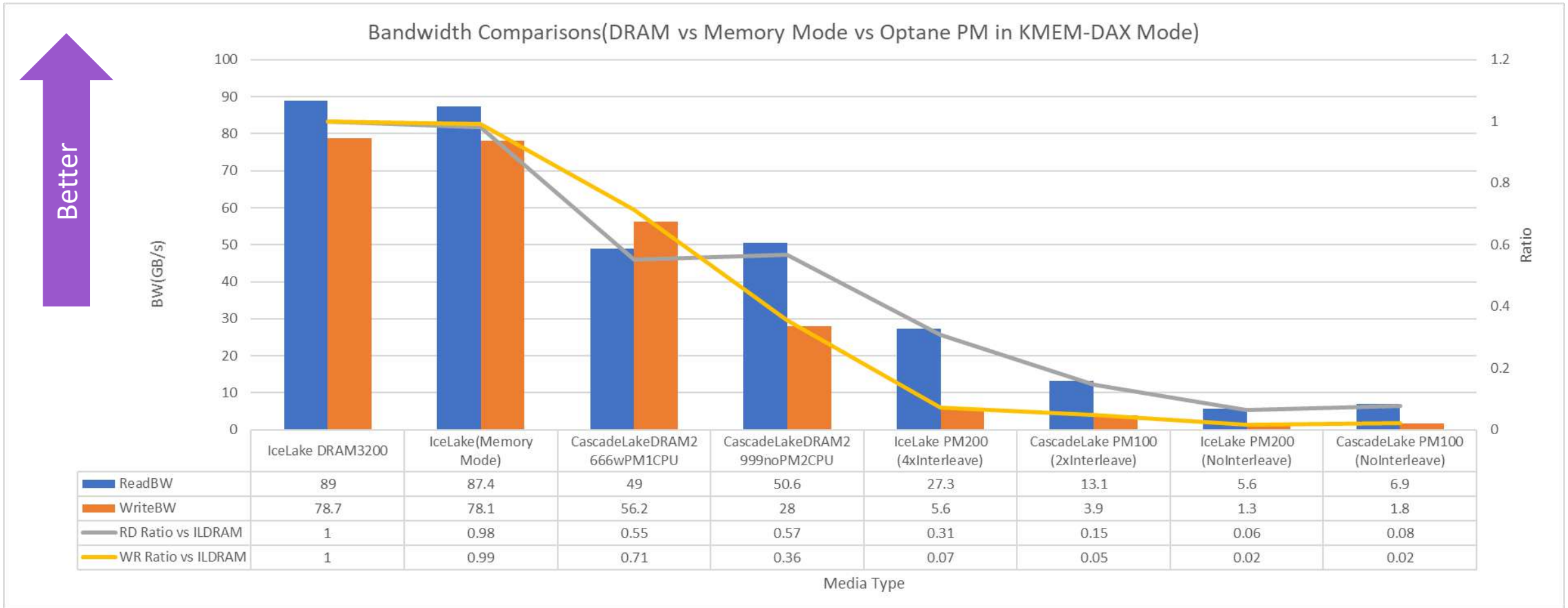
# KMEM DAX

- PMEM を**揮発性メモリ**として活用するための機能
- PMEM のメモリ空間を DRAM とは別の NUMA ノードとして見せる
  - アプリを修正することなく DRAM と PMEM を使い分け可能
- Linux kernel version 5.1 からサポート

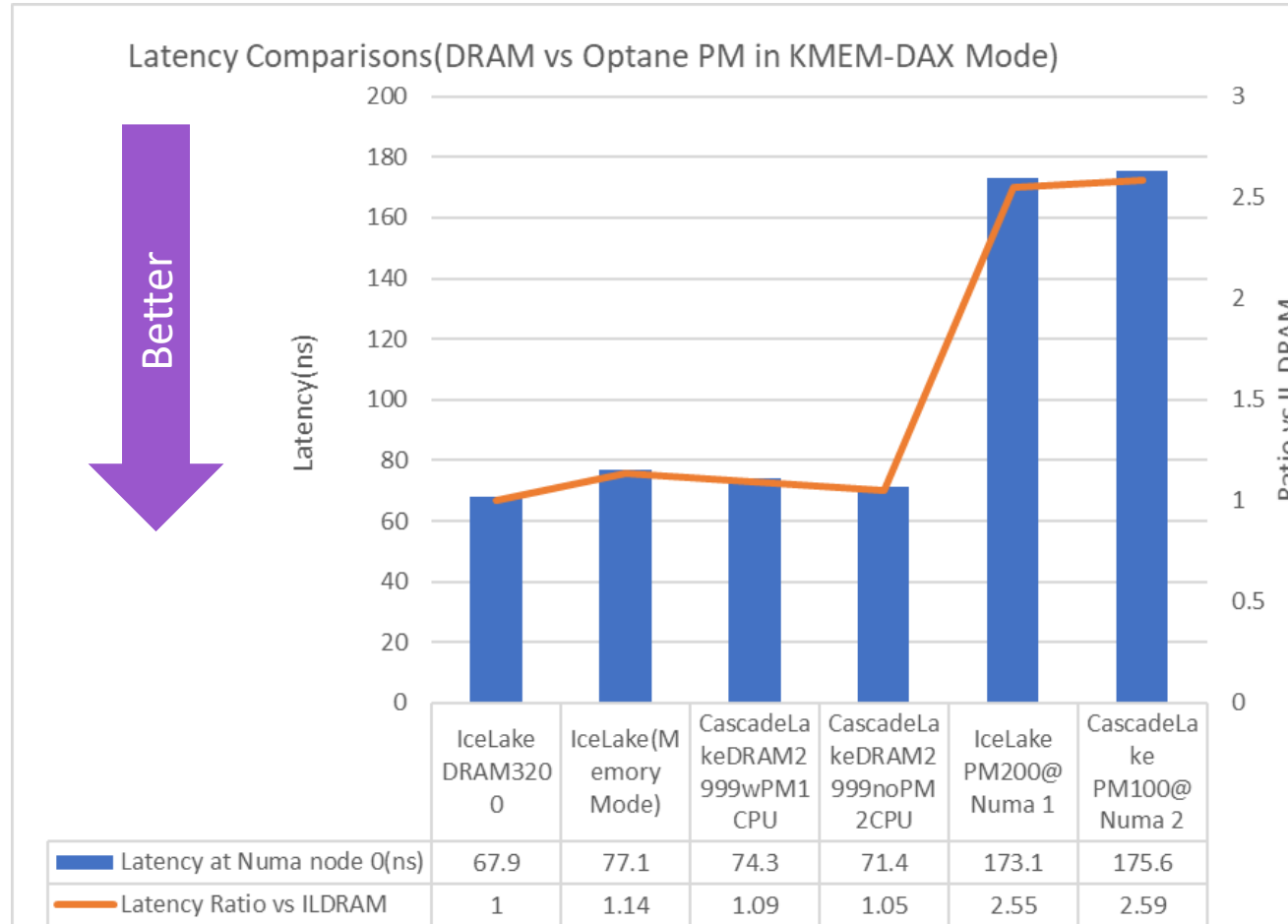
```
[Ice Lake]# numactl -H
available: 2 nodes (0-1)
node 0 cpus: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18
19 20 21 22 23 24 25 26 27 28 29 30 31
node 0 size: 63930 MB
node 0 free: 51967 MB
node 1 cpus:
node 1 size: 507904 MB
node 1 free: 507904 MB
node distances:
node  0  1
  0: 10 17
  1: 17 10
```



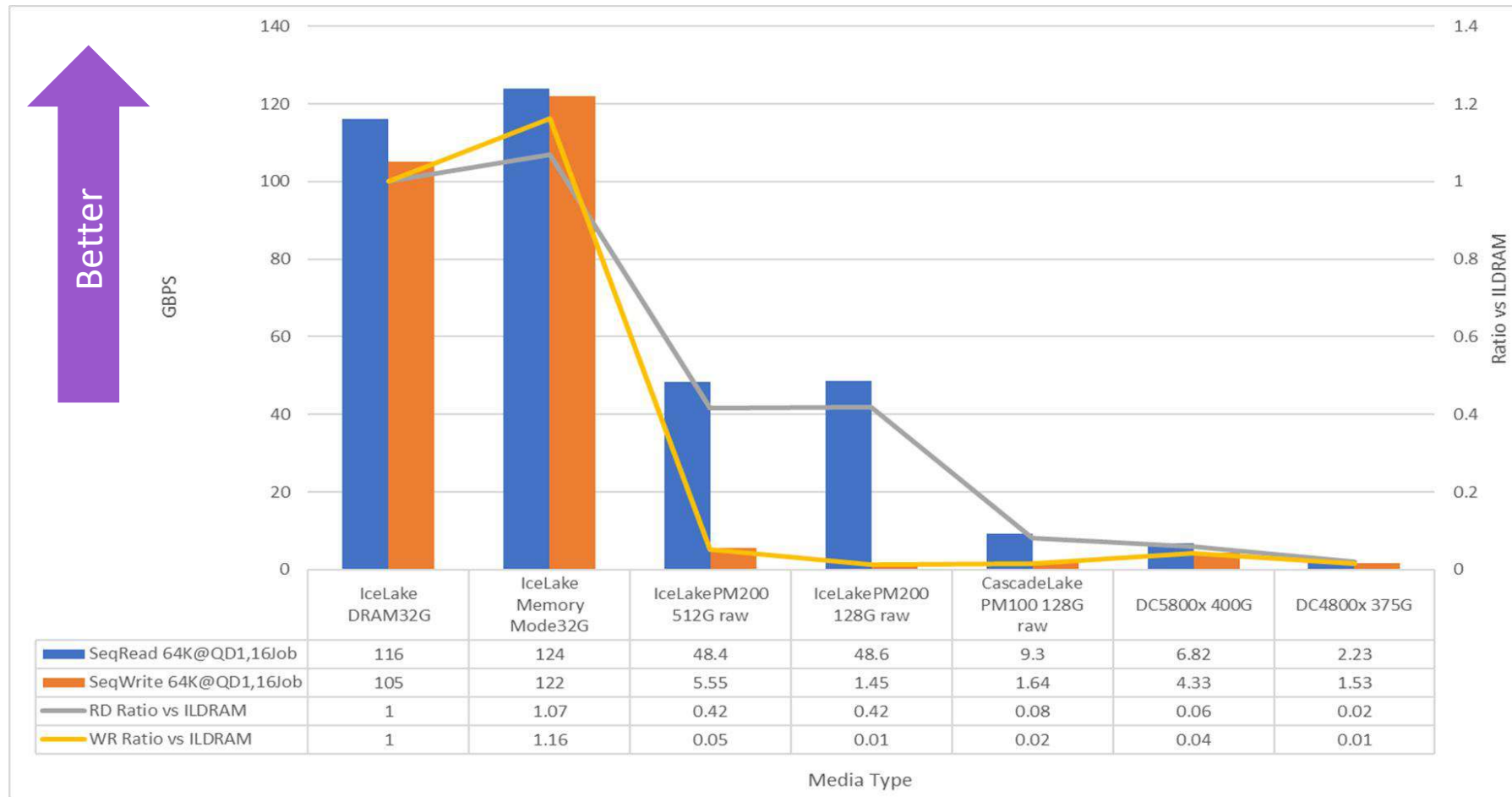
# メモリメモリとしてのバンド幅



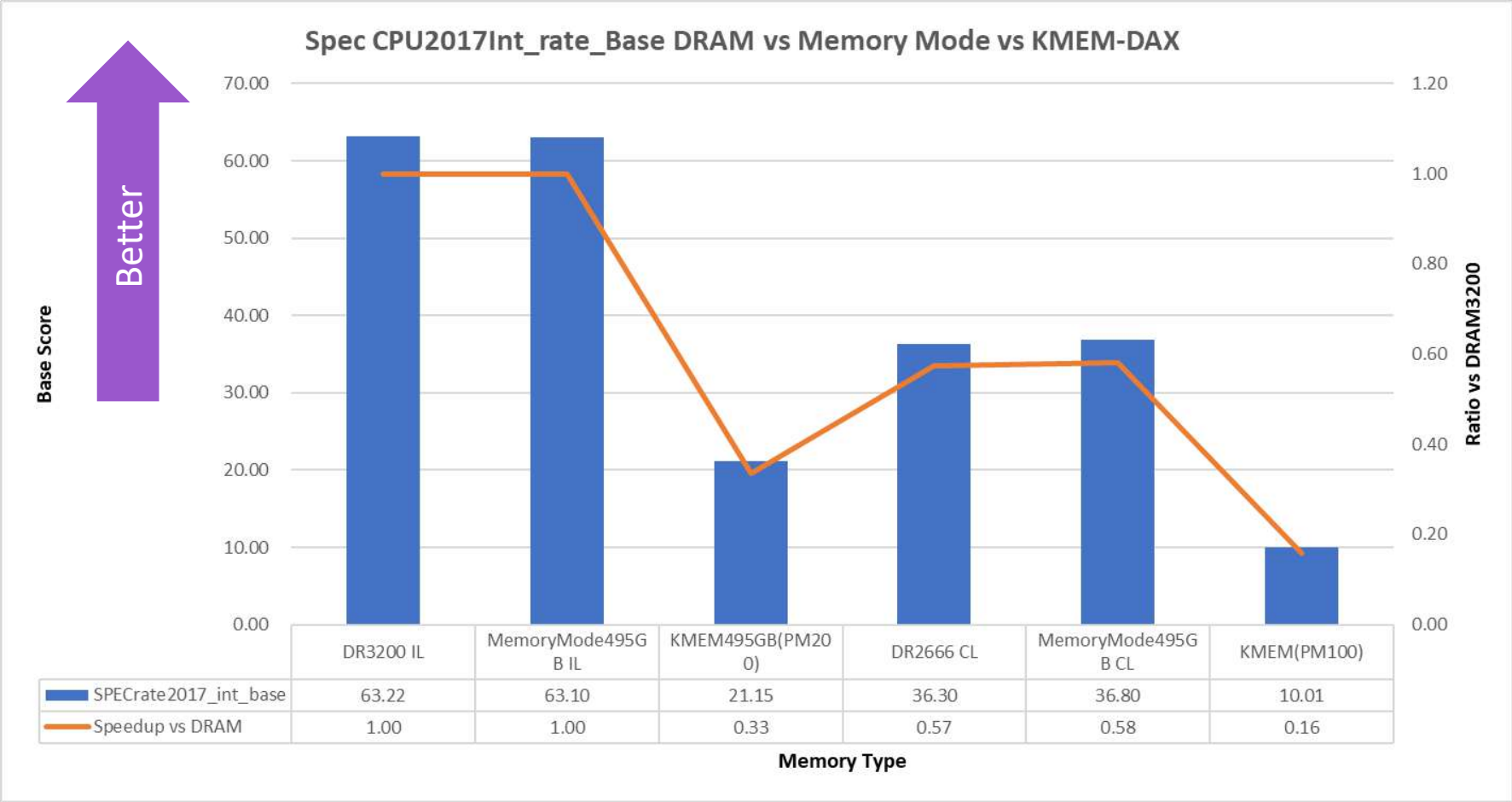
# メインメモリとしてのアクセスレイテンシ



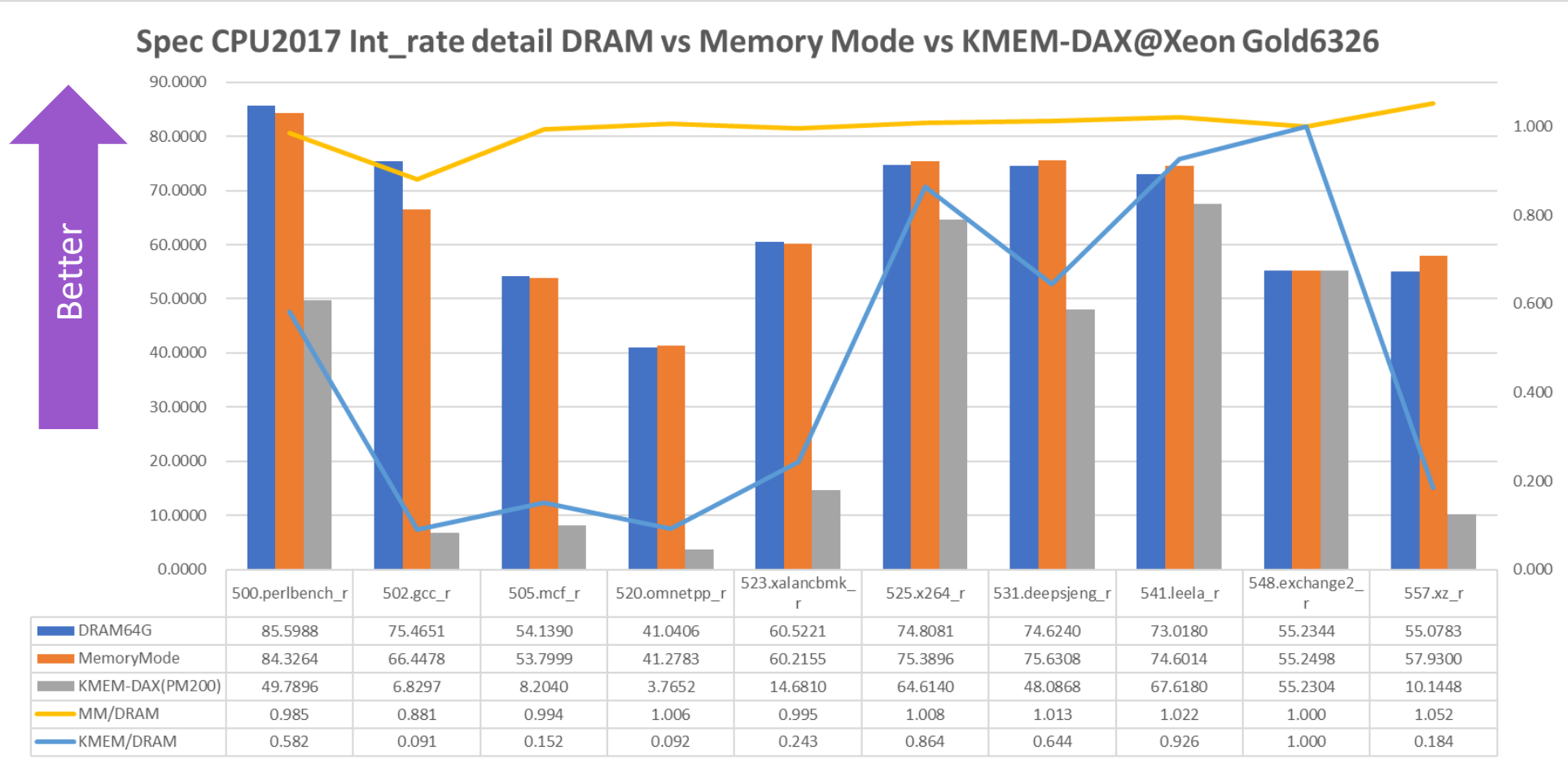
# ブロックデバイスとしてのバンド幅



# SPEC CPU® 2017 ベンチマーク int\_rate

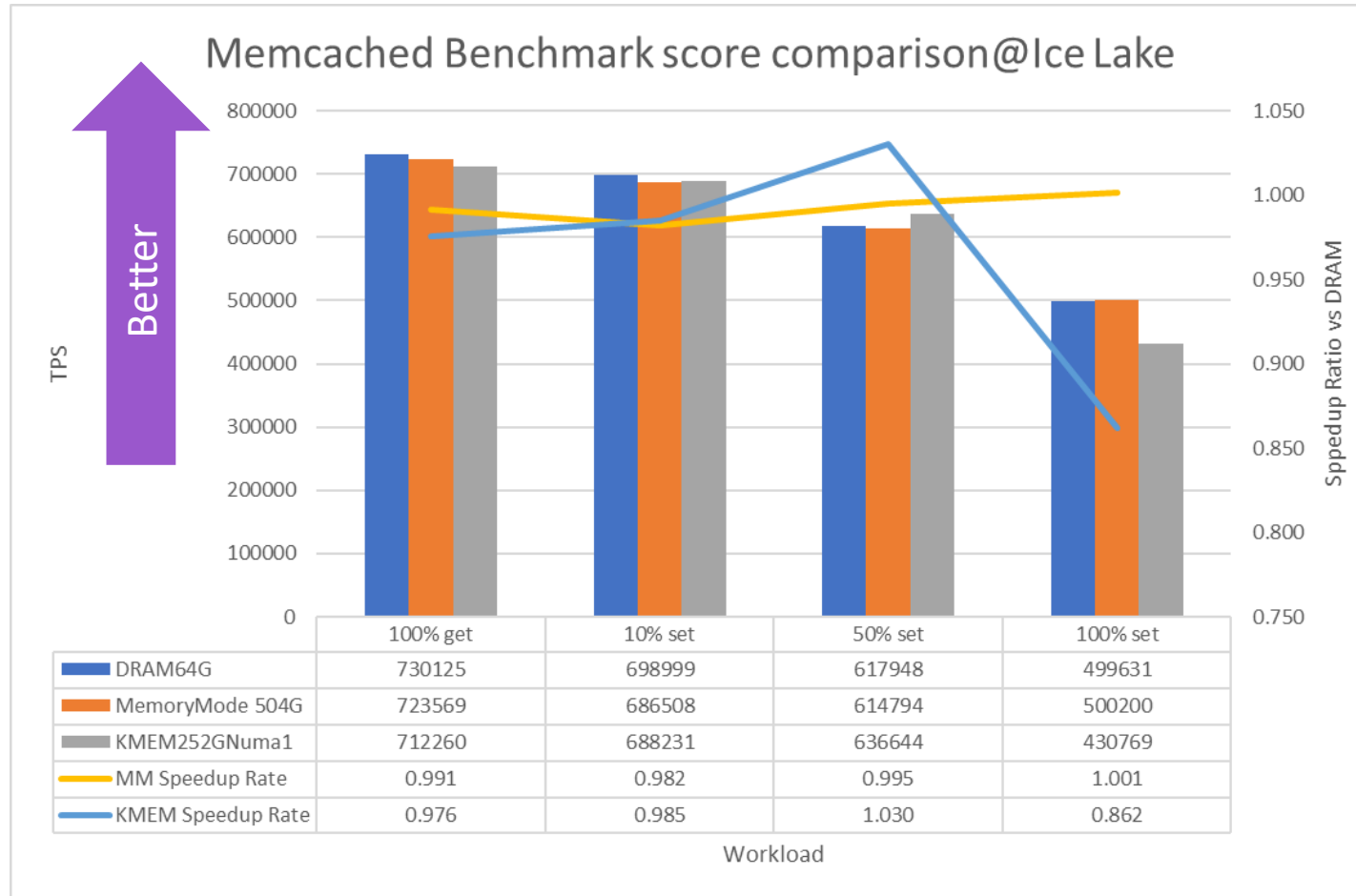


# SPEC CPU® 2017 int\_rate (ベンチマーク毎の詳細)

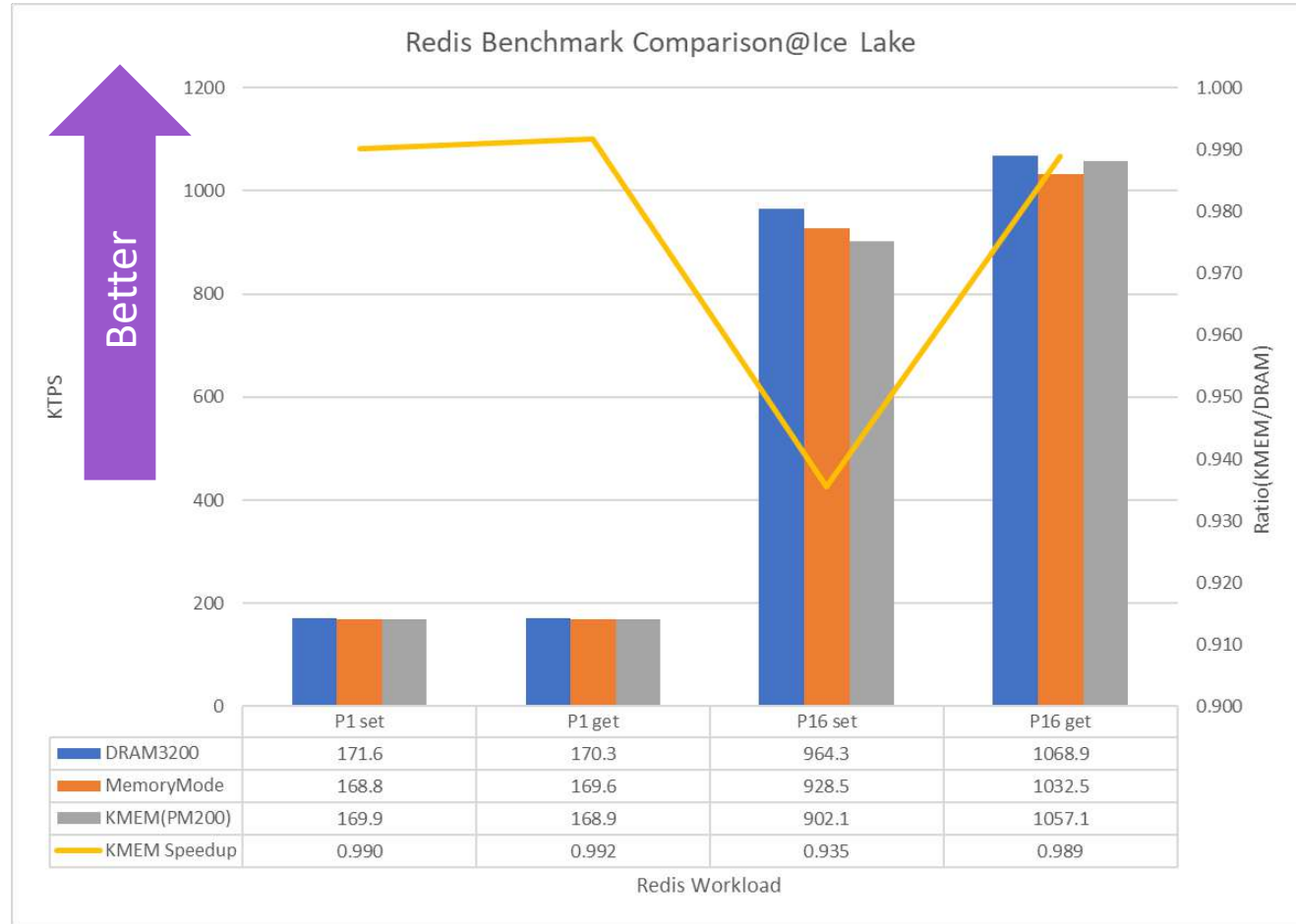




# Memcached (インメモリキーバリューストア)



# Redis (インメモリキーバリューストア)



# まとめ

- PMEM とは「バイトアクセス可能な DIMM タイプの不揮発性メモリデバイス」
- 100 series から 200 series へ性能面・機能面で進化している
- メインメモリ、ブロックストレージ、不揮発性メインメモリとして活用可能
- 基礎性能は DRAM に劣る
- ただし、実アプリに対しては DRAM との性能差は意外と見えない  
⇒ 大容量・低バイト単価という面は DRAM に対してアドバンテージ
  - 実際に仮想化基盤への導入事例なども出てきた ([ソフトバンク](#)、[NEC](#))



# Q&A