

Storage Systems in Preferred Networks

SNIA-J 最新技術動向講演会
2024/10/22

Kota Uenishi

A decorative graphic at the bottom of the slide features a network diagram with nodes and connecting lines overlaid on a faint city skyline. The network diagram consists of numerous small blue dots (nodes) connected by thin white lines, forming a complex web-like structure. The city skyline in the background is rendered in a light blue, semi-transparent style, showing various skyscrapers and buildings.

Preferred Networks 会社概要

ミッション：現実世界を計算可能にする

設立	2014年3月26日
本社	東京都千代田区
代表取締役	西川徹（最高経営責任者）、岡野原大輔（最高研究責任者）
従業員数	約300名（2023年2月）
事業内容	深層学習やロボティクスなどの先端技術を応用したソフトウェア・ハードウェア・ネットワーク技術の研究・開発・販売
主要子会社	株式会社Preferred Computational Chemistry（2021年6月） 株式会社Preferred Robotics（2021年11月） 株式会社Preferred Elements（2023年11月）
出資企業	トヨタ自動車株式会社 ファナック株式会社 日本電信電話株式会社 ENEOS ホールディングス株式会社 中外製薬株式会社 株式会社博報堂DYホールディングス 株式会社日立製作所 三井物産株式会社 みずほ銀行株式会社 東京エレクトロン株式会社



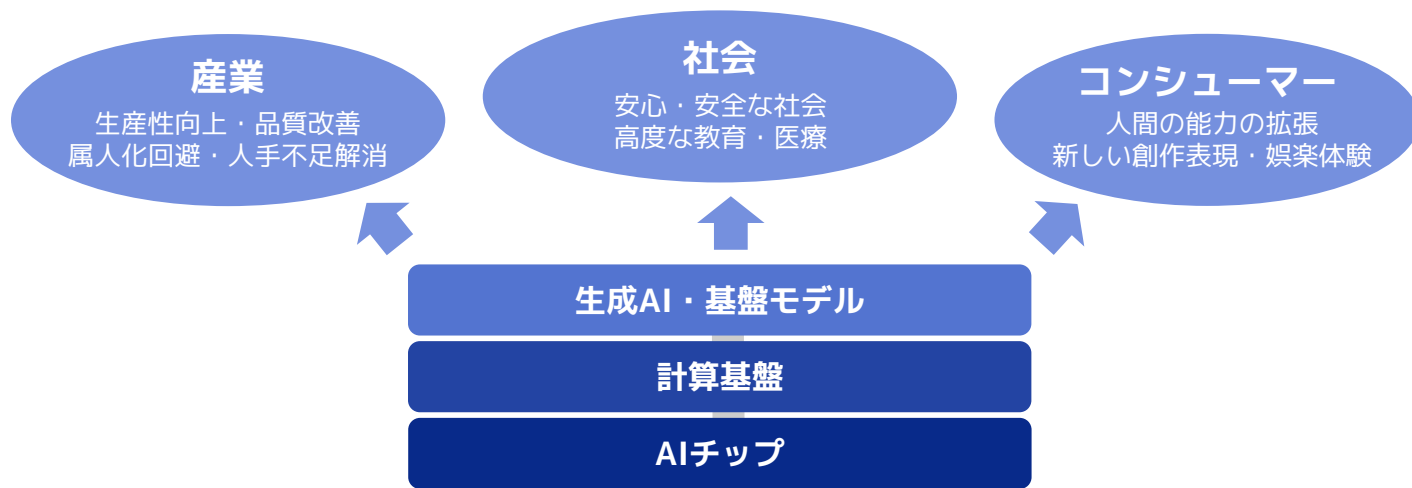
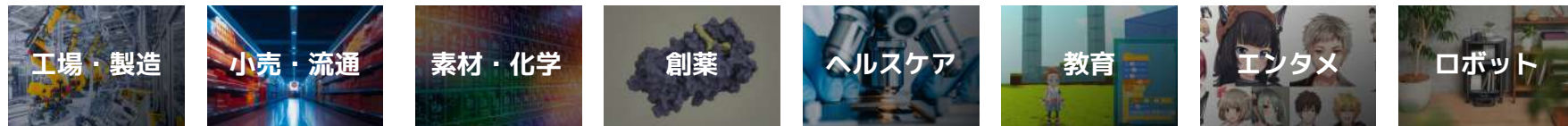
PFNの事業: AI技術のバリューチェーンを垂直統合

PFNは、チップ、計算基盤、生成AI・基盤モデル、ソリューション・製品まで、AI技術のバリューチェーンを垂直統合し、ソフトウェアとハードウェアを高度に融合することで、競争力の高い技術の開発および産業応用を進めています。



PFNの事業: AI技術の水平展開

PFNは、AI技術のバリューチェーンを垂直統合し、産業、コンシューマー、社会に向けて様々な領域でソリューション・製品を水平展開しています。



材料探索を高速化する汎用原子レベルシミュレーター



MATLANTIS

数字で見るMatlantis™



シミュレーション速度

20M 倍

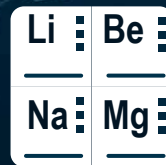
3,000原子で従来のDFT
計算と比較した場合



最大対応原子数

20k 原子

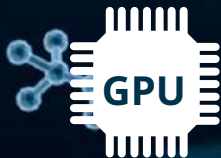
Professional Planで
Pt bulkの場合



対応元素数

72 元素

地殻上に存在する
元素の質量の99.9969%



1 GPU換算の計算量

1,650 年

3,300万以上の構造からなる訓練
データの生成にかかった計算量を
1台のGPUで処理した場合の年数



クライアント数

80+ 組織

企業および学術機関
などの団体



累計原子数

2.6T 原子

2023年1月~3月の3か月間で
シミュレーションした原子の総数

2024年1月現在

Matlantisのニューラルネットワークポテンシャル「PFPP」は、PFNのスーパーコンピュータおよび国立研究開発法人産業技術総合研究所のAI橋渡しクラウド（ABC1）を用いて開発されました。

© Preferred Computational Chemistry, Inc.

エンターテインメント PFN 3D Scan

2022年 日経優秀製品・サービス賞「最優秀賞」受賞

高品質・多様な物品に対応する3Dスキャンサービス

- 2022年6月リリース
- 従来の3Dスキャンでは3Dモデル化が難しかった、透明・金属・黒色の物品などにも対応可能
- すでに3万点近い物品のスキャンを実施

<https://pfn3d.com/>

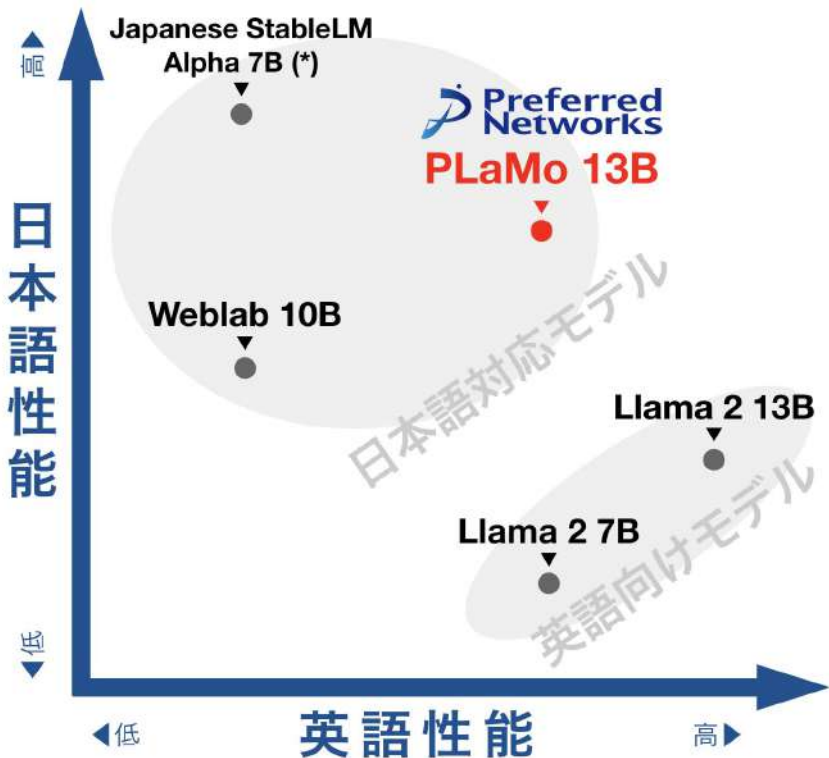


<https://youtu.be/RX4j6wxWev8>



PLaMo-13B PFNのマルチモーダル基盤モデル

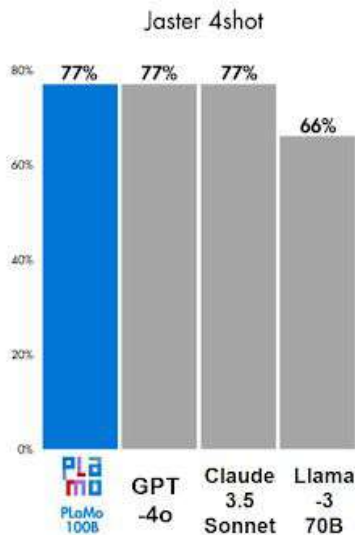
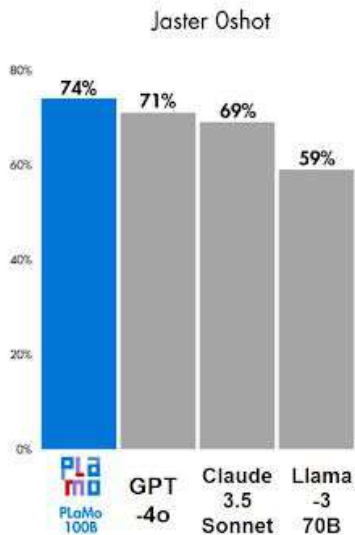
研究/商用利用可能な130億パラメータの大規模言語モデルを去年9月に公開



- 2023年9月公開当時、日英2言語をあわせた能力で世界トップレベル
- 経産省ABCIのスーパーコンピュータ A100 480GPUで1ヶ月弱利用。日本語、英語の1.4兆トークン（数兆文字）を使って学習
- 学習時の実行効率は41%で世界の他の学習基盤と比べても高い（効率的に学習資源を使える）
- 開発実績を元に基盤モデル開発・提供を行う Preferred Elementsを2023年11月に設立



PLaMo-100B-Pretrained の公開



- 研究利用、商用利用可能
- 1000億パラメータ
- 10/15 公開
 - <https://huggingface.co/pfnet/plamo-100b>
- 100% 子会社Preferred ElementsがGENIACプロジェクトの支援のもとフルスクラッチで開発
- 法人向け PLaMo Prime および軽量版 PLaMo Lite を開発中



PFN の計算基盤

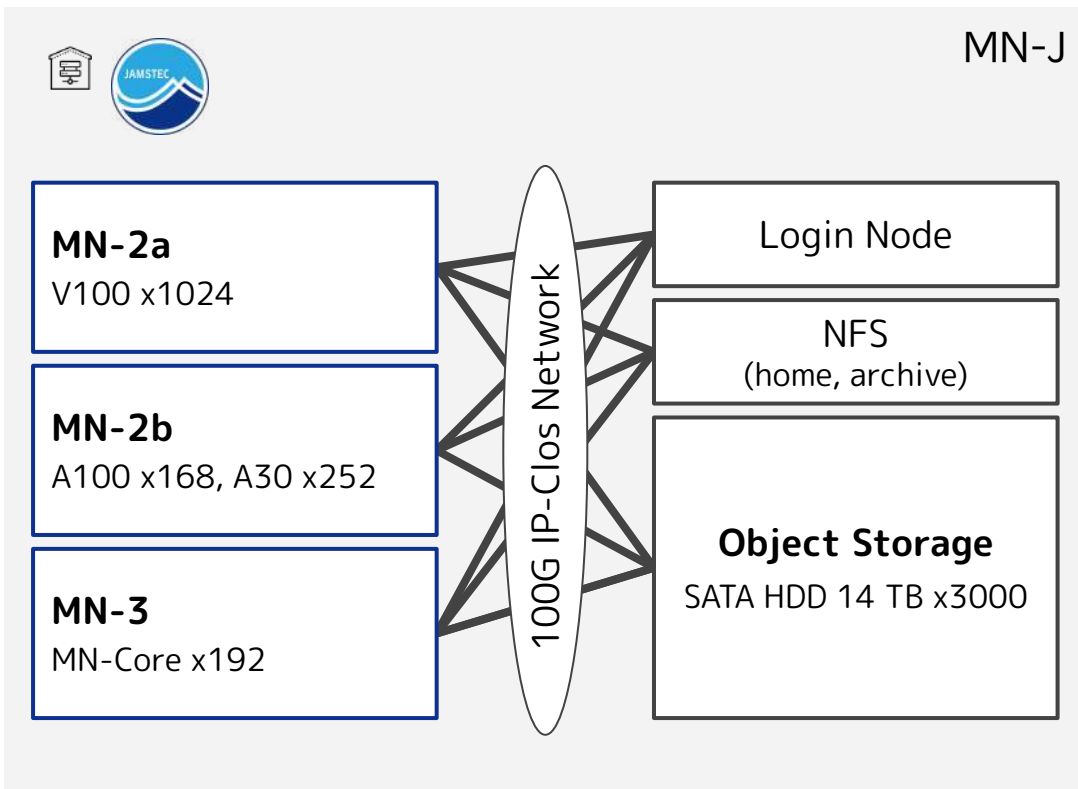


MN-2a, MN-2b

The **GREEN 500** 世界1位!!
(ISC20, ISC21, SC21)



MN-3



パブリック
クラウド



北海道DC
H100 x512






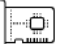
ストレージサーバ

MN-J の計算機







MN-2a

128 nodes
(1024 GPUs)

-  36 CPU Cores
-  DDR4 384 GB
-  V100 (16 / 32 G)
SXM2 x 8
-  100 GbE x 4
RoCEv2
with SR-IOV





MN-3

48 nodes
(192 MN-Cores)

-  48 CPU Cores
-  DDR4 384 GB
-  MN-Core x 4
-  100 GbE x 2
MN-Core
DirectConnect





MN-2b (A100)

42 nodes
(168 GPUs)

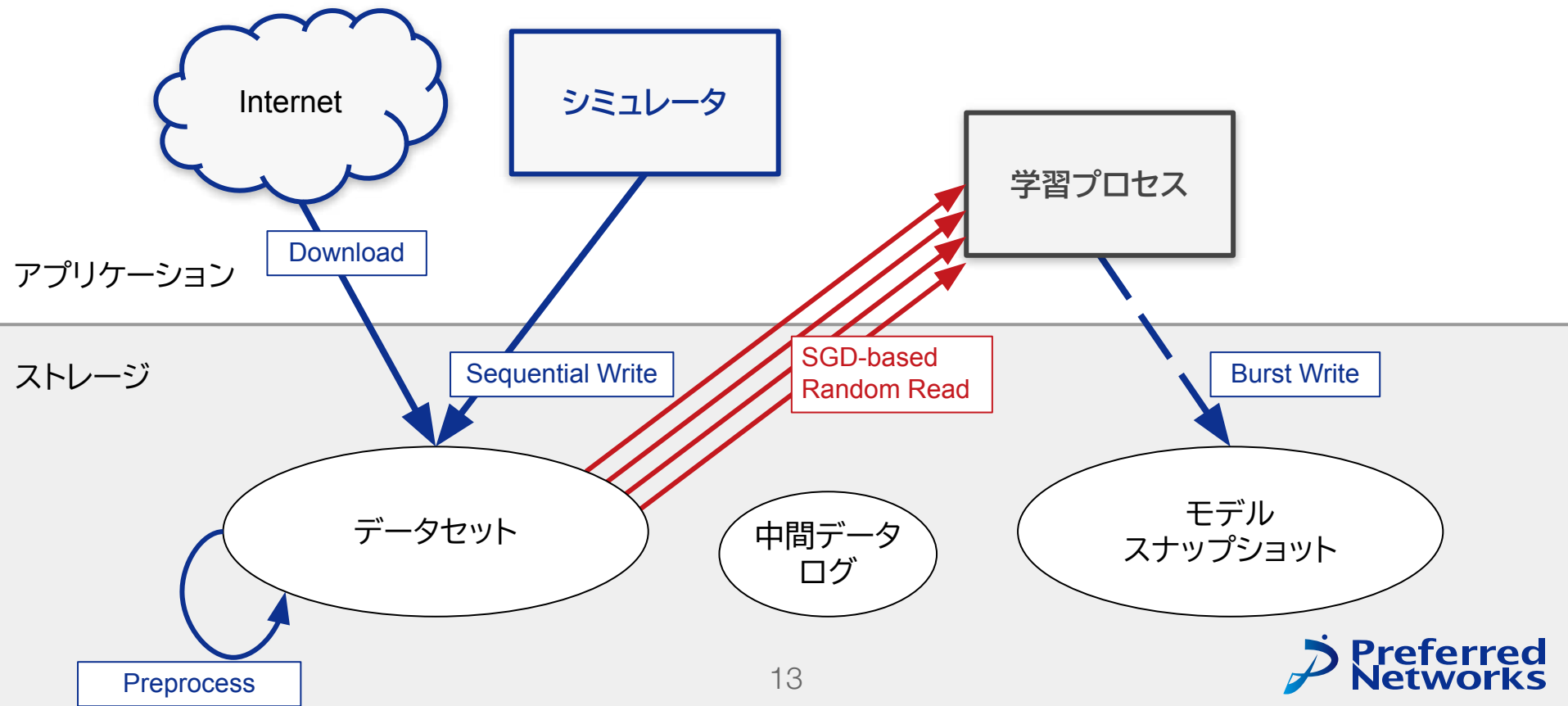
-  128 CPU Cores
-  DDR4 1024 GB
-  A100 (80 G)
SXM4 x 4
-  100 GbE x 2
RoCEv2
with SR-IOV

MN-2b (A30)

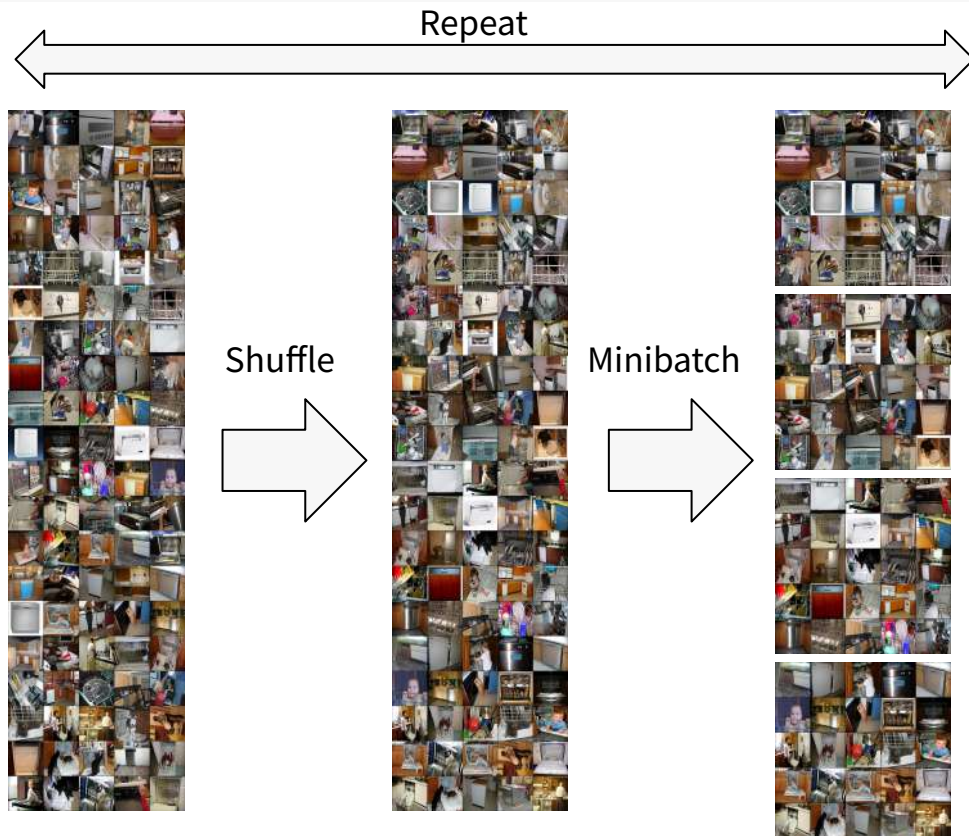
42 nodes
(252 GPUs)

-  80 CPU Cores
-  DDR4 512 GB
-  A30 (24 G)
PCIe x 6
-  100 GbE x 2
RoCEv2
with SR-IOV

ストレージから見た深層学習のワークロード (1/2)



SGD-based Random Read



MNIST

- 60k images
- Typical batch size: 10~100

CIFAR10

- 60k images
- Typical batch size: 64

ImageNet-1k

- 1580k images
- Typical batch size: 32

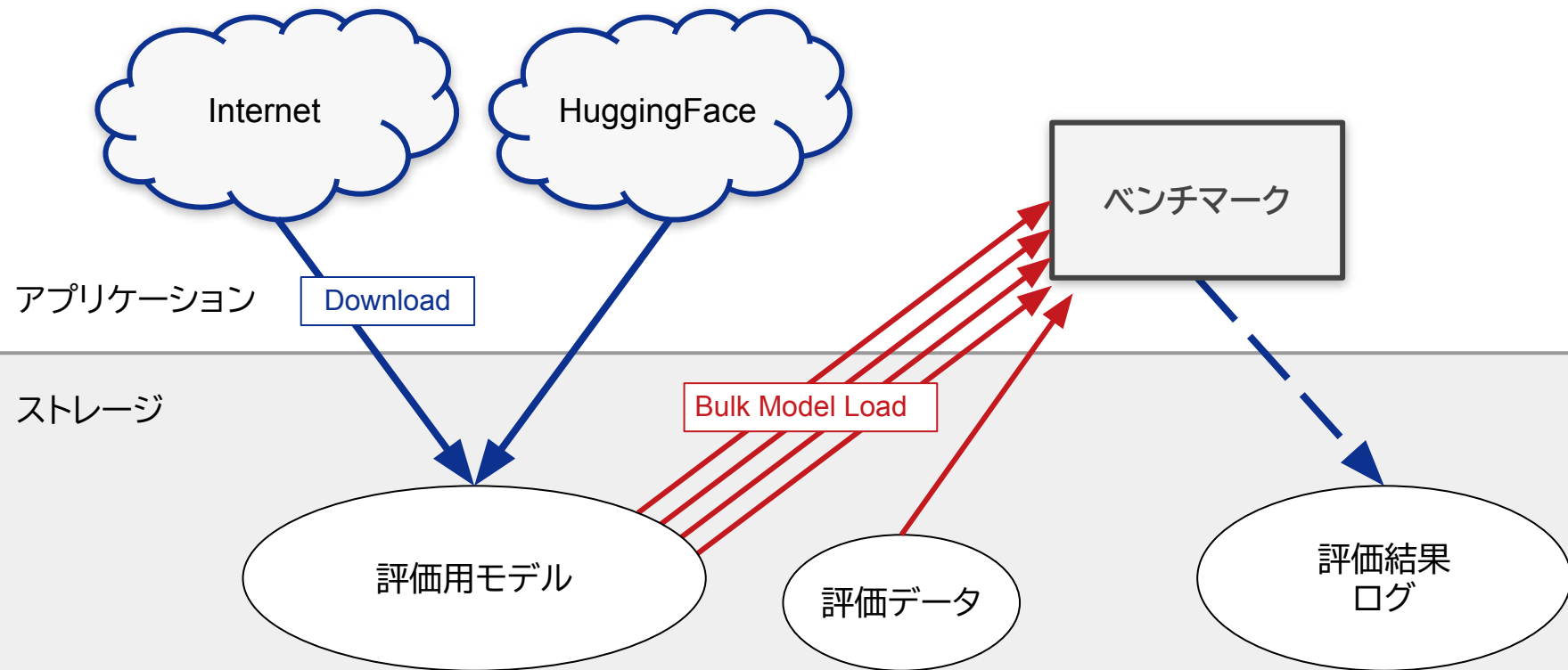
Open Images Dataset (OID)

- 8M images
- Typical batch size: ???

Burst Write

- 故障に備えて学習の途中状態をスナップショットとして保存する
- サイズ
 - CNN-based 画像モデルだと <1GB
 - LLM の場合はパラメータ数にもよるが一般的にはパラメータ数 10B で 120 GB, 100B で 1.2TB
 - 1param あたり 12B
- 速度
 - スナップショットの保存間隔をストレージや計算の MTBF 以内に収めないといけない
 - 一般的には数十分～数時間以内
 - Kubernetes によるプリエンプションもありうる

ストレージから見た深層学習のワークロード (2/2)



モデル評価のワークロード

- PLaMo-100B-Pretrained を始めとする自社モデルの事後学習評価
 - 事後学習のタスクごとに多様なパラメータで何度もモデルをストレージからロードして利用
- HF上にある多種多様なモデルを評価比較
 - 数百GBのLLMを100種類ダウンロードして1000件のタスクを評価する
-

AI インフラに必要なその他の要件

- データセットの保管
 - 画像や動画、テキストベースのデータセット
 - シミュレーション結果
 - 2024/05 時点で総量は 10 PB 程度
- 学習時の高速なデータ読み出し
- 推論時の高速なモデルロード
- システムの水平拡張性
 - 事業が続く限りデータは増え続ける
 - 異なる世代のハードウェアが同一システム下で同居できること
- 単一の名前空間
 - アプリケーションを移植しなくても利用継続できること

要件に対する実装

SGD-based Random Read /

Bulk Model Load

- 低レイテンシ、高い帯域
- 同じデータセットをシャッフルして繰り返し読む
- 大きなモデルを高速ロード



分散キャッシュ
NVMe SSD x 10²



Dataset Generation by Simulation

- 増え続けるデータに対応する
- 永続性、バックアップ

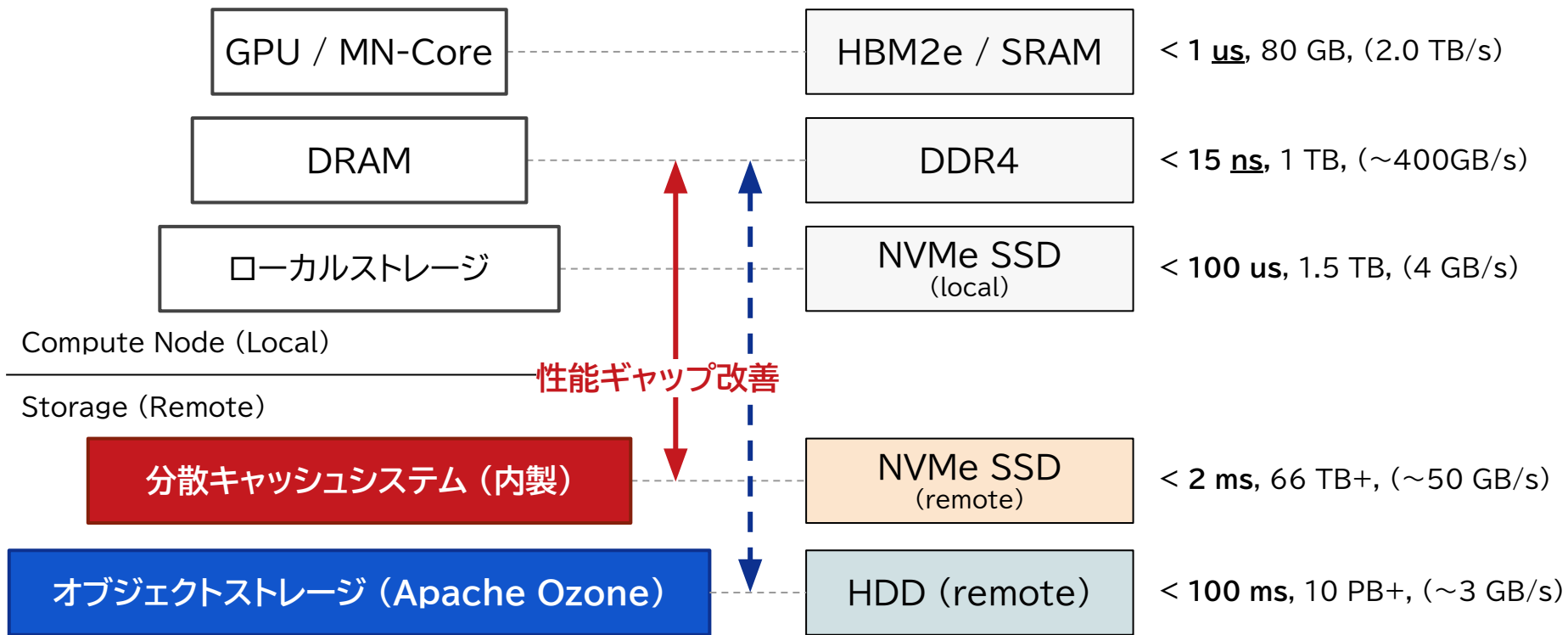


オブジェクトストレージ

Apache Ozone
HDD x 10³



I/O 階層



分散キャッシュシステム (1/2)

- キャッシュシステムの要件
 - SSD の性能を引き出せる
 - データノードを追加して簡単に容量を拡張できる
 - 複数の利用者間で資源を平等に共有できる
- 特徴
 - sendfile(2) によるゼロコピー実装
 - Consistent Hashing によるデータ分散、Envoy による L7 Routing
 - LRU によるデータの入れ替え、クォータ機能
 - Non-POSIX Stateless API (HTTP)

詳細は下記の[スライド](#)・[ブログ](#)を参照



<https://tech.preferred.jp/ja/blog/distributed-cache-for-deep-learning/>
<https://speakerdeck.com/pfn/k8s-tokyo-60-distributed-cache-system>

分散キャッシュシステム (2/2)

- 分散キャッシュシステムはローカルファイルシステム (i.e. JBoF) を想定した設計
 - RAID なし、NVMe-oF なし、永続性なし
 - 高速 (OS からみて ~10us) で大容量な一時領域
 - まずは簡単に HTTP で動くものをつくった
 - NVDIMM が狙っていた領域？
- Discussion
 - NVMe-oF をつかわなかった
 - 運用性
 - Consistent Hashing で負荷分散をしたかったので N:M になる
 - 計算機との RDMA (GPUDirect Storage etc.) はしていない
 - Kubernetes のエコシステム (Envoy) に実装されていない
 - Consistent Hashing で負荷分散をしたかった
 - クライアントは User-land ライブラリにせざるをえなかった
 - FUSE という選択肢
 - コンテナ上の CAP_SYS_ADMIN 問題

c.f. [pfnet-research/meta-fuse-csi-plugin](https://pfnet-research.github.io/meta-fuse-csi-plugin/)

オブジェクトストレージ

- オープンソースソフトウェア選定の検討事項
 - 問題発生時にコードを読める、直せる、データを復旧できる
 - ユーザ空間で動作する
 - 手弁当な小体制でも運用できること
- これまでの運用実績
 - HDFS (~ 2021) → Apache Ozone (2021 ~) → ???

- Cloudera Manager
による省力運用
- 安定性

- Small File Problem
への対処
- S3 API サポート

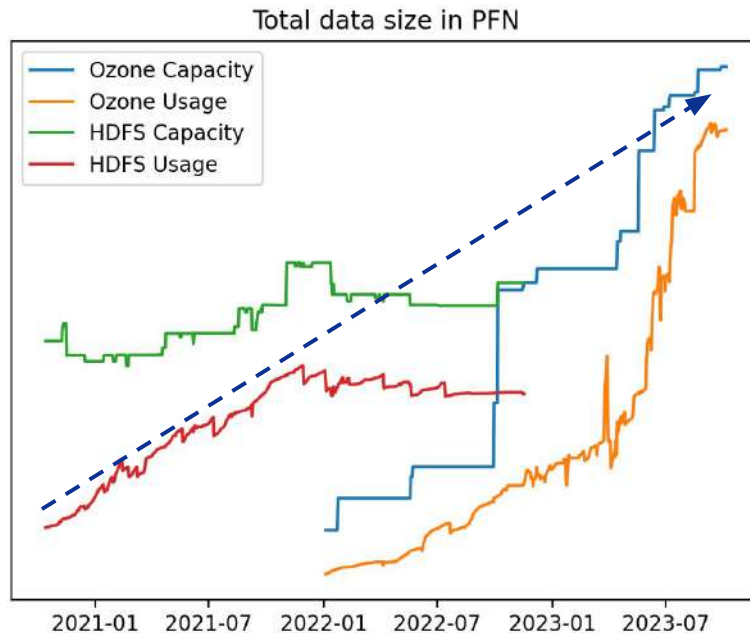
オブジェクトストレージの構成

現在は 3 系統の Apache Ozone クラスタを自社運用中

	クラスタ構成 #1	クラスタ構成 #2	クラスタ構成 #3
総物理容量	18.1 PB	12.1 PB	9.0 PB
データノード台数	36	88	24
デバイス	SATA HDD	SATA HDD	SATA HDD
インターコネク	100 GbE x2	100 GbE x2	100 GbE x2
ソフトウェア	Apache Ozone 1.3.x	Apache Ozone 1.3.x	Apache Ozone 1.4.0
運用開始	2021	2022	2024
構成		OM/SCM HA	OM/SCM HA

社内の容量需要の増加

- ストレージの容量需要は一定ラインで継続
- 毎日データノードの HDD が数本埋まる (10 ~ 20TB/day)



HDD 増設の歴史



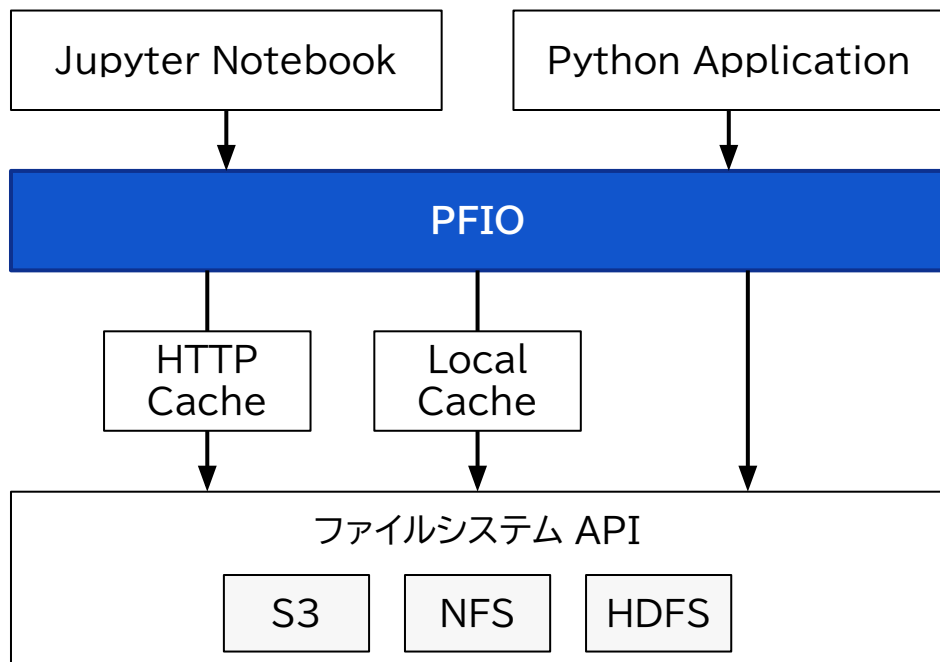
2022



2023

増設はほぼ年 1 回
サーバー数十台、HDD 数百台

インタフェースの抽象化



<https://github.com/pfnet/pfio>

<https://pfio.readthedocs.io/en/latest/index.html>

```
import pfio
```

```
backend_fs = pfio.v2.from_url('ozone://bucket/');  
with pfio.v2.HTTPCachedFS(url = 'http://cache/',  
                           fs = backend_fs) as fs:  
    with fs.open('bin-1m.dat', 'rb') as fp:  
        buf = fp.read()  
        print(len(buf))
```

```
[ozone]
```

```
scheme = s3
```

```
endpoint = ...
```

図: 分散キャッシュと S3 API の抽象化

```
% python3 ./example.py  
1048576  
0.12042808532714844
```



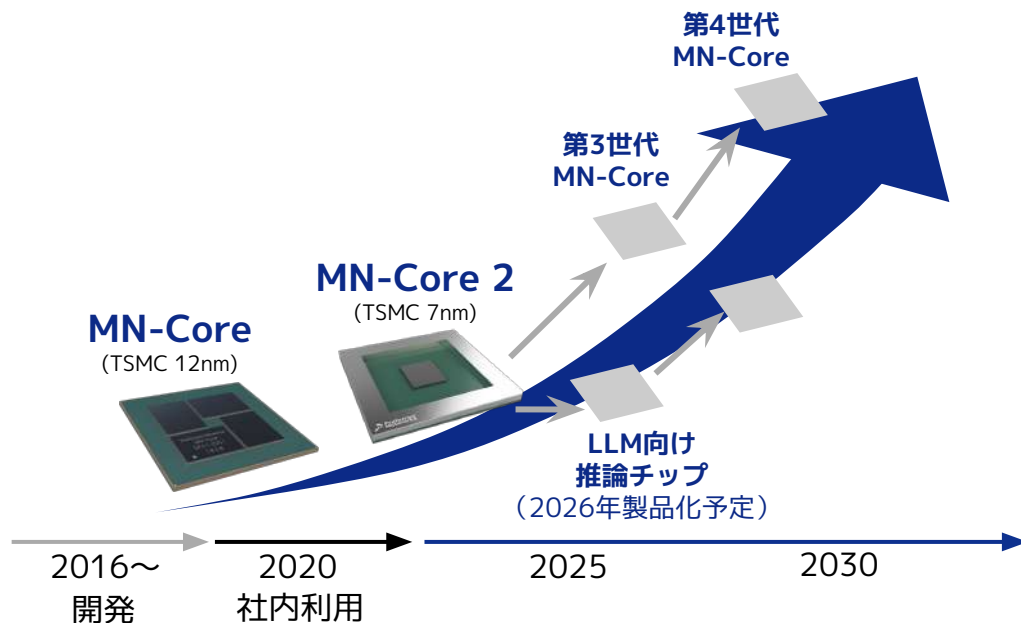
```
% python3 ./example.py  
1048576  
0.017198801040649414
```

図: キャッシュによるレスポンスタイム改善

今後の展開(会社)

MN-Core™ シリーズの進化と今後の展開

- MN-Core : 2020年にスーパーコンピュータMN-3に搭載。社内の研究・開発に利用。
2023年にはMN-Coreによる計算力を試験的に外部ユーザへ提供開始。
- MN-Core 2 : 2023年に試験運用開始。
2024年からMN-Core 2 の販売を開始。クラウドでの提供も予定。



MN-Core クラウドサービス化



The screenshot shows the Preferred Networks website's news section. The main headline reads: "PFN、深層学習用プロセッサMN-Coreを Matlantisの深層学習モデル (PFP) の計算基盤として実装". Below the headline, it states "実際のワークロードで安定して3倍程度の高速化を確認" and "2023.10.16". The article text describes the partnership between PFN and Matlantis, highlighting the performance improvements achieved with MN-Core hardware for deep learning workloads.

<https://www.preferred.jp/ja/news/pr20231016/>



これがMN-Coreね！！

Translate post



12:52 PM · Oct 20, 2023 · 2,304 Views



5



25



1



MN-Core2 クラウドサービス構想



PFN、深層学習を高速化するプロセッサMN-Core 2の開発および、MN-Coreシリーズのクラウドサービス構想を発表

2023年春からPFNのパートナー企業向けにMN-3の計算資源を提供し、順次設備とユーザーを拡大予定

2022.12.14

株式会社Preferred Networks（本社：東京都千代田区、代表取締役 最高経営責任者：西川徹、プリファードネットワークス、以下、PFN）は、深層学習を高速化するディープラーニング・プロセッサ MN-Core™ 2（エムエヌ・コア・ツー）を、東京ビッグサイトで開催されている SEMICON Japan 2022のキーノート講演において本日発表しました。

山田てるみ
@telmin_orca

これがMN-Core2ね！！

[Translate post](#)



12:52 PM · Oct 20, 2023 · 8,982 Views



[Retweet](#) 20

[Like](#) 72

[Bookmark](#) 1



誰もが MN-Core™ シリーズを利用できる AI クラウドサービス



Preferred Computing Platform

Preferred Computing Platform (以下、PFCP) は株式会社 Preferred Networks (以下、PFN) が構築運用する深層学習・AI ワークロード向けのクラウドサービスです。PFNが開発する独自アクセラレータであるMN-Core™ シリーズを唯一利用でき、最先端の性能と効率性を備えています。

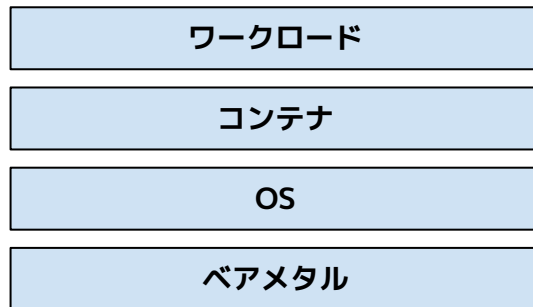
強力な計算ノード

MN-Core ボードを4基搭載した MN-Core サーバを複数専有して利用できます。すべてのノードは深層学習に最適化された高速なネットワークで相互に接続されています。次世代の MN-Core 2 や NVIDIA GPU を搭載したノードも順次提供予定です。

MN-Core 2 サーバの構成	
アクセラレータ	MN-Core 2 ボード × 8
CPU	Intel Xeon SPR 8480+ 56C x 2
Memory	DDR5 64GB x 16 (2TB)
Interconnect	NVIDIA ConnectX-6 (100GbE) × 2

フルマネージドサービス

深層学習・AI ワークロード向けに拡張された Kubernetes クラスタをマルチテナントで利用できます。大規模分散学習から推論サーバの高可用性運用まで幅広く行なえます。ワークロードの状況を観測するためのマネージドなモニタリングサービスも付随しています。



今後の展望(ストレージ)

今後の展望

- SGD-based Random Readワークロード
 - 分散キャッシュの改善
 - 最新世代NVMe SSDの導入
 - RDMA等の導入
- Simulation Data
 - 自社SDSの開発(今はどこも自社開発している)
 - 高密度HDDサーバーの導入
 - HDDはよりアーカイブストレージ化していくが、NANDの単価次第
- LLM Snapshot
 - Transparent HSM-like なシステムの検討

今後の展望: PCIe の高速化

- OSまではソフトウェアが追いついているが、その上のインフラソフトウェアが Gen5 の高速化にすら追従できていない
 - e.g. HDD でいうところの RAID (狭義) 等の技術が追いついておらず、運用しやすい形でファイルシステムAPIを提供する手段が限られている
 - e.g. Gen5 NVMe が安定動作するサーバー機材が限られている (?)
 - e.g. HDD を大量に並べても IO スループットが出ないケースがある
- 400 Gbps の NIC ベンダーは実質 1 社独占であり、SSD だけが高速でも使いにくい状態 (PFNは200Gbpsをスキップしました)
 - 他社が 400 を発表する前に 800 Gbps NIC が発表された
- ストレージ API としては以下を継続利用
 - Relaxed POSIX (NFS, etc)
 - S3 (HTTP)



Making the real world computable