2025年10月2日

『IOWN APN』を活用した疑似遠隔環境における GPU・ストレージ間接続性能テストの詳細と結果

GMO インターネット株式会社 NTT 東日本株式会社 NTT 西日本株式会社 株式会社 QTnet

1. 事前検証の概要

1.1 背景と目的

近年の生成 AI や大規模言語モデル(LLM)の普及により、AI 開発基盤への需要が急激に拡大している。 従来、AI の演算装置(GPU)と大容量ストレージは物理的な隣接配置が必須とされてきたが、データセンター内の設置スペース制約や、演算結果を自社施設に保管したいという多様なニーズに対応するため、地理的制約を超えた分散型 AI インフラの実現が求められている。

本実証実験では、NTT が開発する次世代通信基盤「IOWN(Innovative Optical and Wireless Network) APN(All-Photonics Network)」の高速大容量かつ低遅延性を活用し、GPU とストレージ間の遠隔利用における技術的実現可能性の事前検証として、東京-福岡間の距離を想定した疑似遠隔環境でGMO GPU クラウドの性能テストを実施した。

1.2 各社の役割

GMO インターネット株式会社	GMO GPU クラウドの GPU、およびストレージの提供	
	アプリケーション実装	
NTT 東日本株式会社	IOWN APN 技術提供および実証回線の提供 (※)	
NTT 西日本株式会社	IOWN APN 技術提供および実証回線の提供 (※)	
株式会社 QTnet	データセンター内の実証環境の提供(福岡県福岡市)	

(※) 2025年11月から12月に実施予定の本実証時に使用

1.3 検証スケジュール

事前検証:疑似遠隔環境での性能評価(2025年7月実施) 本実証:実拠点間での接続検証(2025年11-12月予定)

2. 事前検証環境·構成

2.1 物理構成

場所:QTnet データセンター(福岡県福岡市)

GPU: NVIDIA HGX H100ストレージ: DDN AI400X2

ネットワークスイッチ: Arista 7050SX3-48YC8

• 遅延調整装置:OTN Anywhere (東京-福岡間等の距離相当遅延を挿入)

© 2025 GMO Internet,Inc., NTT EAST, Inc., NTT WEST, Inc.,QTnet,Inc. All rights reserved.

Confidential – Citation permitted with acknowledgment. Redistribution prohibited.

2.2 ネットワーク構成

● 接続帯域:100GbE

• 挿入遅延: OTN Anywhere を使用し、東京-福岡間相当等の遅延(0-40 ミリ秒)を疑似的に生成

2.3 疑似遠隔環境の構築手法

物理的に遠隔地にサーバーを設置する代わりに、福岡県福岡市内のデータセンター内に遅延調整装置「OTN Anywhere」を設置し、東京-福岡間等の物理的距離に相当する通信遅延を疑似的に挿入することで、バーチャルな遠隔環境を構築。

3. 事前検証シナリオ

3.1 テストワークロード

本実証では、AI 開発における代表的な2つのタスクを実行し、遠隔ストレージ利用時の性能影響を評価。

3.1.1 画像分類タスク: MLPerf® Training Round 4.0 ResNet (※1 以下 ResNet と表記)

ベンチマーク: ResNet (Residual Neural Network)

• 特徴: ImageNet データセット(約128万枚の学習用イメージを内包)の読み込みと処理を実行

評価指標:目標精度に達するまでの学習時間

3.1.2 大規模言語モデル処理タスク: MLPerf® Training Round 4.1 Llama2 70B (※2 以下

Llama2 と表記)

• ベンチマーク: Llama (Large Language Model Meta AI) 2 70B

特徴: Llama2 70B モデル本体(約 130GB)に対する学習を実行

評価指標:目標精度に達するまでの学習時間

3.2 遅延条件の設定

以下の遅延条件(往復)でそれぞれのタスクを実行し、性能への影響を測定

1. ベースライン: 0 ミリ秒 (隣接配置相当)

3. 長距離遠隔: 20 シリ秒 (東京-沖縄間 約 1400km 相当)

4. 超長距離遠隔:30 ミリ秒 (東京-台北間 約 2100km 相当)

5. 極長距離遠隔: 40 ミリ秒 (東京-マニラ間 約 2700km 相当)

4. 実験結果

4.1 ResNet 画像分類タスクの結果

遅延条件	ベンチマークスコア (分) ^(※1)
0ms	13.80分
15ms	15.55分
20ms	15.95 分
30ms	17.52分
40ms	19.09分

Result not verified by MLCommons Association.

© 2025 GMO Internet,Inc., NTT EAST, Inc., NTT WEST, Inc.,QTnet,Inc. All rights reserved.

Confidential – Citation permitted with acknowledgment. Redistribution prohibited.

4.2 Llama 大規模言語処理タスクの結果

遅延条件	ベンチマークスコア(分) ^(※2)
0ms	24.87分
15ms	24.94 分
20ms	24.95分
30ms	25.01分
40ms	25.07分

Result not verified by MLCommons Association.

5. 考察·分析

5.1 性能影響分析

画像分類タスク、大規模言語モデル処理タスクのいずれにおいても遅延条件を増加させるに従ってベンチマークスコアが悪化していく傾向が観測され、遠隔ストレージ利用時の影響を疑似的に測定、観測することができた。

一方、ベンチマークスコアに与える影響の度合いについては画像分類タスク、大規模言語モデル処理タスクの間で大きな差が存在する。ベンチマークの実装や処理内容を詳細に分析したところ、この差は以下のようにそれぞれのベンチマーク特性、特にベンチマーク内における遠隔ストレージへの I/O 頻度に依存しているものと考えられ、遅延条件を受けやすい(受けにくい)処理が明確に存在することが推察される。

5.2 ResNet 画像分類タスク

ベンチマークが測定開始されてから GPU メモリヘ ImageNet データセットの読み込みを行っており、遅延条件 (ms)増加に伴うストレージからの転送性能の変化がベンチマークスコアの低下に表れている。一方、対象となるデータセットは一般的な AI 学習手法でも用いられる前処理によって生データ (約 128 万枚の学習用イメージ) を学習に適した単一のファイル形式へと整形しているため、その低下度合いは東京・福岡間(15ms)想定で 12% ほどであった。

5.3 Llama 大規模言語モデル処理タスク

ベンチマークが測定開始される以前に GPU メモリへ大規模言語モデル等の読み込みが完了している。測定開始後は主に GPU 上の演算によって完結する処理が多く、ストレージへの I/O は ResNet 画像分類タスクに比べて少量であることからベンチマークスコアの低下度合いも極めて少なかった。

5.4 まとめ

本検証では、画像分類タスクおよび大規模言語モデル処理タスクを対象として、疑似的に遅延条件を付与することにより、遠隔ストレージ利用時の性能影響を測定した。

その結果、いずれのタスクにおいても遅延条件の増加に伴いベンチマークスコアの変化が確認され、遠隔ストレージ由来の遅延影響を観測可能であることを示した。意図したような遅延条件の影響が観測できたこと、加えて本検証で設定した東京-福岡間相当の遅延条件における性能低下は12%程度である。NVIDIA HGX H100 は、前世代の A100 GPU と比較して AI トレーニング性能で最大4倍の高速化が可能とされており(**3)、性能向上率を勘案すると学習における12%の性能低下は許容されると判断した。以上より、本検証を継続する意義が示されたとし、従来回線から IOWN APN への置き換えによる遅延削減効果の検証を目的として、実拠点間でのIOWN APN 接続実証実験の実施を進めていく。

© 2025 GMO Internet,Inc., NTT EAST, Inc., NTT WEST, Inc.,QTnet,Inc. All rights reserved.

Confidential – Citation permitted with acknowledgment. Redistribution prohibited.

6. 今後の展開

6.1 本実証計画(2025年11-12月)

実施内容:東京-福岡間の実際のIOWN APN 回線を使用した検証

比較対象:一般的なイーサネット専用線との性能比較

評価項目:商用実装に向けた実用性評価

6.2 将来的な社会実装ビジョン

本実証の成功により、以下の社会実装が期待される:

- 1. 分散型 AI クラウドの実現:全国規模での AI リソース最適配置
- 2. 災害耐性の向上:分散配置による事業継続性確保
- 3. 新しい社会 NW 基盤の実現:広くあまねく IOWN APN (NTT 東日本・NTT 西日本の「All-Photonics Connect powered by IOWN」) の展開

以上

**1 Unverified MLPerf® Training Round 4.0 Closed Resnet offline. Result not verified by MLCommons Association.

%2 Unverified MLPerf® Training Round 4.1 Closed Llama2 70B offline. Result not verified by MLCommons Association.

The MLPerf name and logo are registered and unregistered trademarks of MLCommons Association in the United States and other countries. All rights reserved. Unauthorized use strictly prohibited. See www.mlcommons.org for more information."

※3 NVIDIA H200 Tensor コア GPU (https://www.nvidia.com/ja-jp/data-center/h100/)