

ホワイトペーパー

FUJITSU AI Zinrai ディープラーニング システム

リファレンスアーキテクチャ テクニカルレポート (ZDLS 200E, ETERNUS A220 編)

Technical Report for FUJITSU AI Zinrai Deep Learning System 200E and ETERNUS A220



はじめに

本書は、富士通 PRIMERGY サーバーをベースとした富士通 AI Zinrai ディープラーニング システム(以降 ZDLS)と、AI ワークフローに最適化された ETERNUS ストレージシステムによるリファレンスデザインについて説明します。これにより ZDLS の構成拡大とともに、AI インフラストラクチャの構築工数削減と導入期間短縮、ONTAP クラウド接続型データストレージによるエンタープライズレベルのデータ管理機能を実現することができます。なお、本書は ZDLS のご利用を前提として作成しています。

本書の読者

本書は ZDLS を導入される社内 SE、関係会社 SE、パートナー様 SE を対象に作成しています。また、ZDLS で採用している Linux OS 及び OSS について習熟していることを想定しており、技術についての詳細な説明は記載していません。

関連ドキュメント

関連するドキュメントとして、以下があります。必要に応じて参照してください。

ドキュメント	概要
ホワイトペーパー Zinrai ディープラーニング システム リファレンスアーキテクチャ デザインガイド (ZDLS 200E, ETERNUS A220 編)	Zinrai ディープラーニング システム リファレンスアーキテクチャ を構築するために必要な操作方法を記載しています。(社内公開)
ZDLS 公開ドキュメント	https://www.fujitsu.com/jp/solutions/business-technology/ai/ai-zinrai/services/deep-learning/#anc-03

- 参考

ONTAP 9 マニュアル <https://docs.netapp.com/ontap-9/index.jsp>

そのほか、各構成部品/OSS のマニュアルを必要に応じて参照してください。

目次

1	本書の概要	1
1.1	リファレンスの目的	1
1.2	リファレンスの利用シーン	1
2	システム構成	2
2.1	全体概要	2
2.2	ハードウェア	2
2.2.1	ZDLS 200E (FUJITSU Server PRIMERGY RX2540 M5)	2
2.2.2	SR-S752TR1.....	2
2.2.3	FUJITSU Storage ETERNUS A220	2
2.3	ソフトウェア.....	3
3	システム構成	4
3.1	ハードウェア構成	5
3.2	ソフトウェア構成	5
4	本リファレンスの特徴・優位性	6
4.1	設計・環境構築が容易	6
4.2	性能面での優位性	8
5	まとめ	8

1 本書の概要

1.1 リファレンスの目的

Zinrai ディープラーニング システム リファレンスアーキテクチャ (以後、本リファレンスデザイン)は、ディープラーニングや機械学習に必要なソフトウェアがプレインストールされた ZDLS プラットフォームを使用したリファレンスデザインです。ZDLS プラットフォームをベースとした PRIMERGY RX2540 サーバーと、ETERNUS ストレージによるマルチノードのディープラーニング環境を、短時間で容易にセットアップを行えることを本レポートで提示しています。

1.2 リファレンスの利用シーン

本リファレンスデザインを利用することで、AI システムのインフラ部分の設計・検証が不要もしくは大幅な削減になります。本リファレンスデザインは、ZDLS をベースに、ストレージやネットワーク機器を拡張したものであり、お客様の要望に応じて不足している部分をカスタマイズして構成することも可能です。

スクラッチから検討・構築・検証を行うより速く確実にインフラストラクチャを構築することができるため、AI によるソリューション開発に、より多くの時間を割くことができます。性能の面においては、最新の NVIDIA Tesla V100S を使用した、ZDLS プラットフォームによるリファレンスデザインであり、ベストオブクラスと言えます。

2 システム構成

2.1 全体概要

本リファレンスデザインは、ZDLS のシステム構成を拡張し、様々な AI ソリューションを搭載できるようにしたオンプレ AI インフラストラクチャのひながたです。ディープラーニングに必要なソフトウェアをインストールしたサーバーとストレージを、ネットワークスイッチを介して 10Gbps ネットワークで接続した、マルチノードのディープラーニング環境の構築例を解説しています。

2.2 ハードウェア

2.2.1 ZDLS 200E (FUJITSU Server PRIMERGY RX2540 M5)

最新のプロセッサ 第二世代 Intel Xeon Scalable Processor Family を搭載した 2U デュアルプロセッササーバーFUJITSU Server PRIMERGYRX2540 M5 は NVIDIA Tesla V100S GPU (HBM 32GB)を最大 2 基搭載可能な、AI や HPC のワークロードに最適な高性能サーバーです。ディープラーニングに必要な OS・ドライバ・ソフトウェアがインストール済みの状態で出荷されるため、導入後すぐにディープラーニング環境を使用できます。

※本リファレンスでは外部ストレージと接続する構成を示していますが、サーバー単体によるスモールスタートで良ければ ZDLS 200E だけで小規模な学習環境を利用開始することができます。



2.2.2 SR-S752TR1

SR-S752TR1 スイッチは、中小規模バックボーンスイッチとして最適な多ポート L2 スイッチです。10 ギガ/マルチギガインタフェースを標準で搭載し広帯域化への対応、サーバーと同じ方向のエアフロー設計（前面吸気背面排気）、電源の冗長化対応など、サーバーやストレージを効率的に収容可能な高機能多ポート L2 スイッチです。



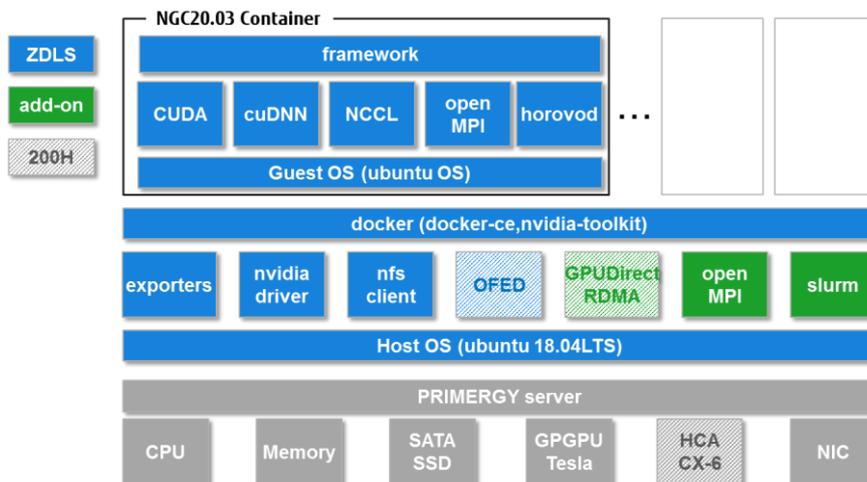
2.2.3 FUJITSU Storage ETERNUS A220

NR1000 A220 は、スモールスタートに最適なエントリーオールフラッシュです。最大容量 2.2 PB の優れた拡張性・RAID4、RAID-DP/TEC 採用による、運用中の 1 ディスクドライブ単位での活性増設・12 Gbit/s 高速ドライブインターフェースをサポートします。



2.3 ソフトウェア

ZDLS はディープラーニングや機械学習に必要なソフトウェアがプレインストールされたプラットフォームです。このため、マルチノードで動作させるための設定を追加するだけで、短期間で本システムの GPU コンピューティングプラットフォームのセットアップが可能です。青で示した ZDLS で構築済みのソフトウェアスタックに、緑で示すソフトウェアを追加します。斜線でハッチングしている部分は ZDLS 200H の場合にさらに追加される部分ですが、本書の説明範囲には含みません。



3 システム構成

本リファレンスデザインでは、NR1000 A220 ストレージシステム (2 コントローラによるクラスタ構成)、ZDLS 200E(PRIMERGY RX2540 M5 サーバー) (2 台)、および SR-S752TR1(1 台)を使用しました。図に示すように、各 PRIMERGY RX2540 M5 は、10GBase-T ネットワークアダプターを介して、計 4 つの 10GbE 接続で SR-S752TR1 スイッチに接続され、NFS ストレージアクセスを行います。NR1000 A220 ストレージシステム内の、2 台の各ストレージコントローラーも、計 4 つの 10GbE リンクを使用してネットワークスイッチに接続します。

ZDLS 200E には、出荷時にインストールされている ZDLS ソフトウェアスタックに加えて、ジョブスケジューラである slurm と、並列計算を行う Open-mpi をサーバーに追加でインストールし、マルチノード環境を構築します。NR1000 A220 ストレージシステムには、SVM(Storage Virtual Machine)を構築し、その中で FlexGroup を作成しています。サーバー当たり合計 10TB のボリュームを割り当てています。

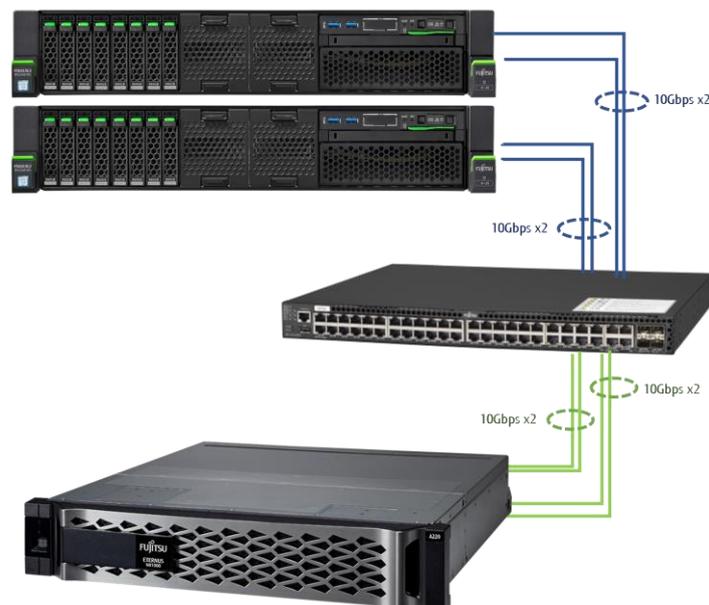


図 1 システム構成

3.1 ハードウェア構成

今回のリファレンスデザインでは以下の機器を使用しています。

表 1 ハードウェア構成

ハードウェア	数	補足	主なスペック
PRIMERGY RX2540 M5	2	CPU: Dual 24 core Intel Xeon Gold 6226R Memory: 18 x DDR4 32GB 2933MHz GPU: 2 x NVIDIA Tesla V100S 32GB Storage: 2 x SATA SSDs Network: 2 x 10Gbase-T Power Consumption: max 1,020[W]	NVIDIA Tesla V100S: GPU クロック:最大 1597 MHz メモリクロック:最大 1107 MHz メモリ帯域:最大 1134 GB/s
NR1000 A220	1	2 コントローラによる HA ペア 12x SSD	4KB ランダム Read:225,300 IOPS 4KB ランダム Write:128,030 IOPS
SR-S752TR1	1		スイッチ容量 248Gbps 最大パケット転送能力:18,452 万 PPS

3.2 ソフトウェア構成

ZDLS プラットフォームおよび、今回のリファレンスデザインで使用される主なソフトウェアは以下になります。

表 2 ソフトウェア構成

	Software	Version
OS	Ubuntu	18.04.4 LTS
NVIDIA 関連	Nvidia-driver	440.64.00
	cuda	10.2
Docker	docker-ce	19.03
Framework	TensorFlow,PyTorch,MXNet	
リファレンスデザイン向け (ZDLS に追加)	Slurm	19.05.4
	openmpi	3.1.5
	munge	0.5.13

4 本リファレンスの特徴・優位性

本リファレンスデザイン特徴とマルチノードのディープラーニング環境を構築する場合の性能の面での優位性について述べます。

4.1 設計・環境構築が容易

ML や DL などの開発を行う AI のシステムでは一般的に GPU を使った処理が中心になります。しかし GPU を使うためのドライバや各種ソフトウェアの選択や構築手順は複雑で、実際の構成 検討や構築にノウハウや多くの時間を要します。通常であれば、技術調査 - 構成品選定 - 構成検討 - 構築手順作成 - 装置構築 - 動作検証といった手順によってシステムが使用可能となりますが、本リファレンスデザインを使用することで、技術的な調査や構成品の選定・構成の検討といった作業は概ね不要となります。必要な作業は、構築手順確認 - 装置構築 - 動作検証となり、AI のシステムのインフラストラクチャ構築の時間を大幅に短縮できます。

ZDLS としてソフトウェアがプリインストールされたサーバーは、Ubuntu OS インストールおよび、DL ソフトウェアのインストール・設定・動作検証が済んだ状態で出荷・納品されます。装置構築・動作検証のフェーズにおいて、一般的に手作業でインストール・設定を行うことと比較して、多くの時間を短縮することが可能です。

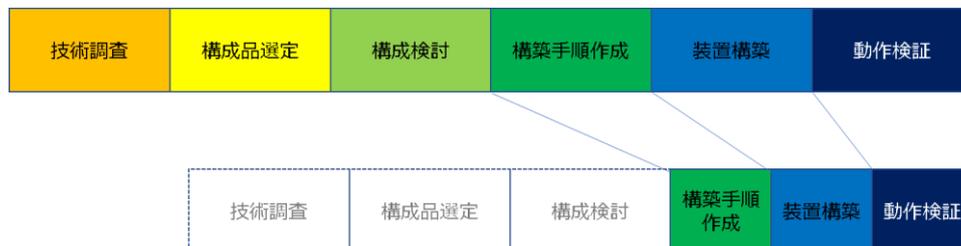


図 2 システム構築完了までの時間の比較

本リファレンスデザインでは、数～数十種程度のコマンド入力により、マルチノードのディープラーニング環境を構築します。各装置に対して行う主な設定は以下になります。

サーバーRX2540 M5 に対しては、ネットワーク等の設定や、7 種のパッケージ・ライブラリの追加インストールを、すべてのサーバー上で行います。Ubuntu など、Linux OS 上でのソフトウェアインストールやネットワーク等の設定の経験がある作業であれば、短時間の作業で、環境構築と動作確認が完了します。

- apt/wget 環境の確認
- ssh root login の設定
- hostname/hosts の設定
- 公開鍵認証の設定
- known_hosts の追加
- 自動更新の抑止
- データ LAN 用のネットワーク設定
- 追加ライブラリのインストールと設定
 - munge

- libmunge-dev
- libmunge2
- libnuma-dev
- ntp
- slurm
- openmpi

ストレージ装置 NR1000 A220 に対しては、サーバーからの NFS アクセスを行うための 30 種程度のコマンド入力、構築が完了します。主な設定内容は以下になります。NR1000 シリーズの操作の経験のある作業者の場合、60 分程度の作業で、環境構築と動作確認が完了します。

- interface group の作成
- network port の作成
- broadcast domain の作成
- aggregate の作成
- SVM の作成
- nfs の設定
- export ポリシーの作成
- logical interface(Data LIF)の作成
- Flex Volume の作成
- Flex Group の作成

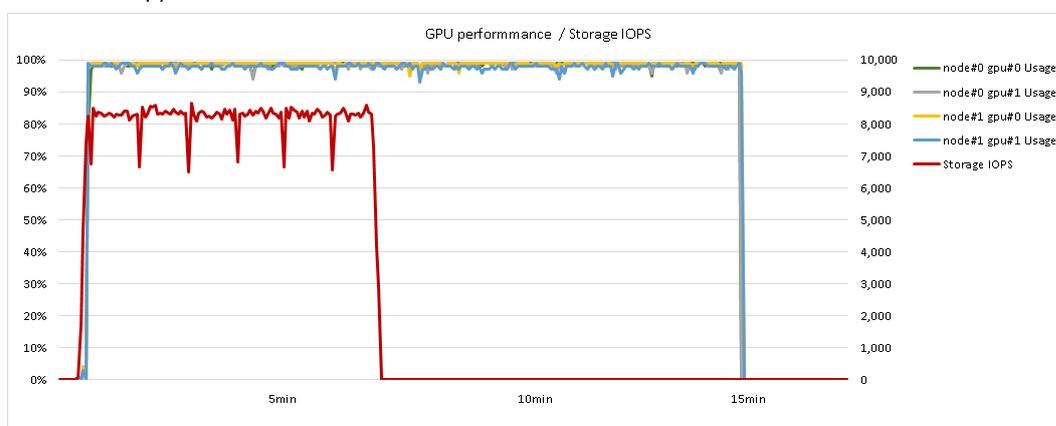
ネットワーク装置 SR-S752TR1 に対しては、接続されるサーバー・ストレージに対するアグリゲーション設定・VLAN 設定を行うことで、構築が完了します。SR-S シリーズの操作の経験のある作業者の場合、30 分程度の作業で、環境構築と動作確認が完了します。

- ポート情報の設定(VLAN)
- リンクアグリゲーションの設定
- LAN 情報の設定

4.2 性能面での優位性

一般的なサーバーシステム同様に、CPU 性能の高いサーバーや、I/O 性能の高いストレージシステムを使用しても、簡単にシステム全体の性能を向上させることができるとは限りません。システム全体の性能を上げるには、CPU や GPU にデータを供給し計算結果を出力するまでの一連の動作があるため、ストレージやネットワークの性能にも注意が必要となります。AI システムにおいて、計算処理を行う際に重要なのは、処理に必要なデータが格納できることと、CPU や GPU の処理に応じてデータを供給できることであるため、ストレージやネットワーク装置の選定と設定が性能に影響すると考えられます。システムを中心となる GPU 搭載サーバーにストレージやネットワークを含めたシステムとして提示している、本リファレンスデザインを活用することで、GPU の性能を十分に引き出すことが可能となり、システム全体の性能を最適化できます。

本リファレンスデザインで構築した 2 ノードのシステム上で学習を行ったときの、ストレージ装置への I/O アクセスと GPU の使用率は以下になりました。GPU 使用率は学習開始後すぐに 100% 近くまで達し、終了までその値を維持しています。本リファレンスデザインの構成において、学習の中心となる GPU の能力を最大限に発揮できており、ストレージアクセスやネットワークアクセスにボトルネックが無いことになります。一方で、もっともストレージ装置への I/O アクセスが発生している、学習の最初 epoch でのストレージ装置の負荷は 8,000 IOPS 程度であり、サーバー数をさらに追加した構成であっても十分に余裕のあることがわかります。途中からストレージ装置の IOPS の値がほぼ 0 になっているのは、一度、学習に必要なデータを読み出すと、次の epoch 以降の学習は、ストレージ装置へのアクセスは不要となるためです。学習のフレームワークは mxnet 20.03-py3、学習モデルは Resnet50 v1.5 を使用しています。



5 まとめ

本レポートでは、ZDLS サーバーをベースとした環境に対し、いくつかのソフトウェアを追加でインストール・設定するだけで、マルチノードのディープラーニングプラットフォームを容易に構築できることを示しました。また、サーバー-ストレージ間の I/O 性能やネットワーク性能がボトルネックになることなく、GPU 性能を十分に発揮できる構成であることも示し、構築面と性能面での優位性を示しています。本レポートでは、GPU に NVIDIA Tesla V100S を使用したモデルについて説明していますが、ZDLS プラットフォームでは NVIDIA Tesla T4 や NVIDIA Tesla V100 SXM2 を使用したラインナップも用意しており、これらに対しても同様の手順で構築が可能です。ベースとなる ZDLS プラットフォームに対して、本リファレンスデザインを使用することで、お客様から求められるさまざまな規模やコストに対応できるディープラーニング環境を提案することが可能になります。

また、今回のリファレンスデザインから得た知見を元に、今後は、Kubernetes・Kubeflow の技術・導入支援を検討していきます。お客様にとって、いかに簡単に使い始められるか・将来の拡張性・安心のサポートといった観点から、引き続き ZDLS を活用したプラットフォームの検討・支援を行っていきます。