Data Pipeline Considerations for Hyperscale AI

Santosh Rao Senior Technical Director, AI & Data Engg

Apr 2019





At distance L₅ ~ 20 m (~ 60 m), the resolution is 5 cm/pixel



rear



AI - Impactful

Rapid AI adoption world wide







\$59.8 Billion Growth of AI software market in 2025







190% AI patents grew by over 5 yrs



71% Automation potential in Manufacturing

4th la startur

4th largest number of Al startups in Berlin







Refining Analytics doesn't lead to Artificial Intelligence





Analytics vs A.I





•

© 2018 NetApp, Inc. All rights reserved. — NETAPP CONFIDENTIAL —





.

 \mathbf{i}

.

. *

· · ·

. .

.



© 2018 NetApp, Inc. All rights reserved. — NETAPP CONFIDENTIAL —

.

.

.

. . .





Source: Kenji Doya Complementary roles of basal ganglia and cerebellum in learning and motor control © 2018 NetApp, Inc. All rights reserved. — NETAPP CONFIDENTIAL —

Cellular Network Traffic Scheduling with Deep Reinforcement Learning

Sandeep Chinchali¹, Pan Hu², Tianshu Chu³, Manu Sharma³, Manu Bansal³, Rakesh Misra³ Marco Pavone⁴ and Sachin Katti^{1,2}

¹ Department of Computer Science, Stanford University

² Department of Electrical Engineering, Stanford University

³ Uhana, Inc.

⁴ Department of Aeronautics and Astronautics, Stanford University {csandeep, panhu, pavone, skatti}@stanford.edu, {tchu, manusharma, manub, rakesh}@uhana.io



Figure 1: Time-variant congestion patterns in Melbourne.

Efficient Large-Scale Fleet Management via Multi-Agent Deep Reinforcement Learning

Kaixiang Lin Michigan State University linkaixi@msu.edu Renyu Zhao, Zhe Xu Didi Chuxing {zhaorenyu,xuzhejesse}@didichuxing. Jiayu Zhou Michigan State University jiayuz@msu.edu

com



Human-level control through deep reinforcement learning

Volodymyr Mnih¹*, Koray Kavukcuoglu¹*, David Silver¹*, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

Curriculum Learning

Yoshua Bengio¹ Jérôme Louradour^{1,2} Ronan Collobert³ Jason Weston³ (1) U. MONTREAL, P.O. Box 6128, MONTREAL, CANADA (2) A2IA SA, 40BIS FABERT, PARIS, FRANCE (3) NEC LABORATORIES AMERICA, 4 INDEPENDENCE WAY, PRINCETON, NJ, USA

When a large language model is trained on a sufficiently large and diverse dataset it is able to perform well across many domains and datasets. GPT-2 zero-shots to state of the art performance on 7 out of 8 tested language modeling datasets. The diversity of tasks the model is able to perform in a zero-shot setting suggests that high-capacity models trained to maximize the likelihood of a sufficiently varied text corpus begin to learn how to perform a surprising amount of tasks without the need for explicit supervision.⁵

Special computation properties



Number Representation



- Deep Learning is Empirically Scaleable (Baidu)
- Computationally Homogenous
- Constant runtime & memory use
- Highly Portable
- Easily Baked into silicon → "Semiconductor Rennaisance"



- Relax Precision : Small integers are better
- Relax Synchronization : data races are better
- Relax Communication : sparse communication is better
- Relax Cache Coherence : incoherence is better
 "Olukutan,Stanford Neurips Keynote 2018"

Data Pipeline for AI Workflow

Scale and Optimize Each Stage of the Data Pipeline



16 NetApp Insight © 2018 NetApp, Inc. All rights reserved. NetApp Confidential - Limited Use Only

Edge to Core to Cloud

Seamless data management



Full Scale-out ONTAP AI with DGX-1

24-node A800 cluster, driving 108 DGX-1's





Full Scale-out ONTAP AI with DGX-2

24-node A800 cluster, driving 36 DGX-2's





Dense Cabinet for Al

Deliver Dense Cabinet for Exascale Al

Challenge:

• Data centers facilities lack the power and cooling for the latest highperformance AI computing infrastructure.

New Capability:

- Dense and Modular Dynamic Density Control (DDC) liquid-air cooled cabinet for Al
- This cabinet combines the efficiency of water with the flexibility of air, cooling up to 52kW of power load in a 45U cabinet.
- Cabinets can be deployed in any environment.
- Provides clean-room environment, guaranteed air flow, integrated security, and fire suppression.



scalematrix & ddc cabinets support 45U & 52kW







NetApp

ARTIFICIAL INTELLIGENCE Trends >

Artificial Intelligence Timeline

Why Now?



Artificial Intelligence



Artificial Intelligence



Artificial Intelligence

Machine Learning

K-Means Logistic Regression Decision Trees Random Forests

Deep Learning

CNN RNN LSTM

GAN

Q Learning Deep Reinforcement TD Learning Learning DQN Prioritized Experience Replay Actor Critic Policy Gradient







AI Requirements Discovery

WHAT PROBLEM ARE YOU SOLVING?

Defining the AI/DL Task

INPUTS		BUSINESS QUESTION	AI/DL TASK	EXAMPLE OUTPUTS		
				HEALTHCARE	RETAIL	FINANCE
		ls "it" <u>present</u> or not?	Detection	Cancer Detection	Targeted ads	Cybersecurity
Text Data	Images	What <u>type</u> of thing is "it"?	Classification	Image Classification	Basket Analysis	Credit Scoring
		To what <u>extent</u> is "it" present?	Segmentation	Tumor Size/Shape Analysis	Build 360° Customer View	Credit Risk Analysis
⊗ .⊗	Ļ	What is the likely <u>outcome</u> ?	Prediction	Survivability Prediction	Sentiment & behavior recognition	Fraud Detection
Video	Audio	What will likely satisfy the objective?	Recommendations	Therapy Recommendation	Recommendation Engine	Algorithmic Trading

SOME KEY DECISIONS TO MAKE

FACTOR	DESCRIPTION			
DL Challenge	Supervised or unsupervised, classification or regression, # of labels?			
Architecture	What is the simplest architecture I can use?			
Training Model	How am I going to tune my neural net? Kinds of non-linearity, loss function and weight initialization? Best training framework?			
Data Quantity	How much data will be sufficient to train my model? How do I go about finding that data and is it evenly balanced?			
Data Quality	Is my data directly relevant to the problem & real world data.			
Data Labels	Is training data is labeled same as raw data sets, how do I 'featurize'?			
Data Similarity	Is data same length vectors or does it require pre-processing?			
Data Storage & Access	Where is it stored, locally and on network Data pipeline? How do I plan to extract, transform and load the data (ETL)?			
Infrastructure	Cloud, On-premise, Hybrid. GPUs, CPUs or both? Single or distributed systems? Integration with languages, ent. apps/ databases.			

Al Requirements Discovery - 101



Al Requirements Discovery - Data Science

- What are your top AI use cases / applications?

- Do you use ML or DL or both?

- Data set footprint and growth pattern?

- Data types used – image, video, text, time-series, ...

- How big is the data set?

- Composition of the data set – small files, large files, or a combination?

- Are your GPUs fully utilized?

- Bottleneck in the workflow?
- Use distributed training?
- If yes, environment setting?

- Foot print of your compute cluster – # of GPUs, CPUs

- Physical memory of the m/c
- Cached copy of data?

- Where is the data stored?

- Where/how will training read the data from?
- Any read/write/latency performance targets?
- Maintain versions of datasets?

- ML/DL framework used?

- Software/ framework to speed up data pre-process?

- Software/ framework to control the cluster / environment / storage?

- Version control of library / API for ML/DL?

- NGC or vendor native?
- Container orchestration at scale?

- Success metrics of training?

- Model selection.
- Model footprint?
- # of concurrent models.
- Model deployment platform for Inference.

- How do you pre-process your data set?

- Time spent to process data?
- Tools & vendors used?

Al Requirements Discovery - Data Engineer



DL Model Training Flow



33 © 2019 NetApp, Inc. All rights reserved.

NetApp

Meet Diverse Needs across Data Science and Infra Functions

Data Scientists

Real-world Data for AI / DL

Need Agile Model DevOps :

- Refreshed access to Production Datasets
- Hybrid Cloud for Model Dev
- Distributed Data Science
- Diverse Data Sources
- Model and Data Parallelism
- Multi-Tenant Model Serving
- From Model to Application

Data Architects

Future Proof Architecture

Seek Extensible Architecture:

- Architecture Scales from PoC to Production
- Future Proof to absorb technology changes
- TCO for Massive Datasets
- Maximize Utilization
- Global Scale

Data/IT Admins

Lowest TCO in face of shrinking budgets

Balance Cost & TTM :

- Leverage vs. Dedicate HW Infra
- Stable Operations & Upgrades
- Supported Components & Ecosystems
- Diverse needs across Big Data, AI/DL and HPC

Balanced Architecture to Deliver for Stakeholders

Deployment Options for Al

Where is the source data?

Move Data into Al Platform	Data In-Place	Co-Lo Solution	Source Data in Cold Storage	Cloud Deployment
 80 - 90% deployments HDFS, Splunk, NoSql , Lustre, GPFS → ONTAP AI (NFS, GPUs) Leverage Data Movers to move data into ONTAP AI FabricPool for data tiering 	 All reside and deploy on ONTAP Concept of Unified Data Lake Data on ONTAP AI FabricPool for data tiering 	 Greater control of data NPS solution Data on NPS, GPUs/ Services on the Cloud 	 Data is moved in from cold data tiers for model training Move data from StorageGrid into ONTAP AI 	 Data / GPUs provisioned on the public clouds Use Cloud Volumes Service for file services GPUs on Cloud for compute



Move Data to AI Platform

Data movement from HDFS, MapR-FS, GPFS, Lustre, S3 to AI Platform

Move Data into ONTAP AI



In-Place Data with Hybrid Cloud Option

Unified data lake serving CPU and GPU Compute Clusters

Data In-Place

AI Platform in Co-Location Data Center

Driven by Power, Cooling, Lease options, As-a-Service based, Consumption based

Co-lo Solution

Cold Data Tier as Data Source

Source Data in Cold Storage

NetApp

Federating ML & DL

Unifying Machine Learning & Deep Learning

In-Place Data Pipeline Federating ML & DL

Software Platforms for Al

H2O.ai with ONTAP AI

NetApp partner for ML platform

H2O.ai is transforming the use of AI with software with its category-creating visionary open source machine learning platform.

Key Features:

43

- Predictive analytics with visualization
- Build More Models in Less time
- Machine Learning: supervised and unsupervised algorithm

H2O.ai and ONTAP AI Joint Solution

- Tested and validated the H2O software running on ONTAP AI
- Eliminates design complexities and guesswork with a solution brief.

Why ONTAP AI for H2O.ai ?

- **Deliver AI at scale:** Scale from zero to 100TB deployments in a matter of seconds with a simple, easy-to-use cloud interface.
- **Reduce Risk:** Instantaneous snapshots & cloning allow data scientists to experiment with datasets without risking data loss.
- **Dynamically change service levels:** Enhance performance or reduce OpEx with the ability to dynamically change storage service levels.
- **Build more models in less time:** Reduce the time it takes to develop accurate, production-ready models by automating time-consuming data science tasks.
- **Train across hybrid cloud :** Seamlessly move H20 Driverless AI workloads between on-premises ONTAP AI infrastructure and Cloud Volumes Service.

Allegro.ai with ONTAP AI

NetApp partner for DL platform

Allegro.ai offers the first true end-to-end Al product life-cycle management solution with a focus on DL applied to computer vision.

Key Features:

44

- Automated annotation
- Experiment management + data management
- Continuous learning

Allegro and ONTAP AI Joint Solution

- Tested and validated the Allegro software on ONTAP AI.
- Eliminates design complexities and guesswork with a validated reference architecture.

Why ONATP AI for Allegro?

- Accelerate time to completion: ONTAP Al's high-throughput AFF and NVIDIA GPUs enable you to train more models in less time while extracting maximum performance from your hardware.
- **Increased GPU utilization:** ONTAP AI helps cut down on GPU idle time by keeping the GPUs engaged more often.
- **Simplified deployment:** ONTAP AI's prescriptive architecture eliminates design complexities and enables independent scaling of compute and storage.
- **Fast and big cache size:** ONTAP AI provides large pools of flash storage that can be used for caching the data further reducing the overall job completion times.

Allegro + NetApp

Optimizing Deep Learning Pipelines At Scale

Deep Learning computer vision & sensor fusion platform

allegro.ai helps data scientists / engineers manage & operationalize the full lifecycle of deep learning - from data-set creation to model post deployment self-learning.

Allegro.ai

- Fully integrated end2end platform
- Data Management / Experiment Management / Resource Management

Allegro Data Abstraction

Multi-Site Real World Optimized DL Pipeline Allegro + NetApp

Edge to Core to Cloud

Seamless data management

AI Public Resources

Collateral

- ONTAP AI Reference architecture NVA-1121-design
- ONTAP AI Deployment guide NVA-1121-deploy
- Building a Data Pipeline for Deep Learning WP-7299
- Edge to Core to Cloud white paper WP-7271
- Al with GPUs on AWS & Cloud Volumes Service TR-4718
- Scalable AI Infrastructure WP-7267
- Designing a Data Pipeline for Your AI Workflows WP-7264
- Solution brief SB-3939
- IDC Technology Spotlight paper
- Cambridge Consultants success story

netapp.com/ai

Recent Blogs

- Your Guide to Everything NetApp at GTC 2019
- Al Across Industries: Manufacturing, Telecom & Healthcare
- How to Configure ONTAP AI in 20 Minutes with Ansible
- Unifying Machine Learning and Deep Learning
- Bridging the CPU and GPU Universes
- Is Your Infrastructure Ready for AI Workflows in Production?
- Accelerate I/O for Your Deep Learning Pipeline
- Addressing AI Data Lifecycle Challenges with Data Fabric
- Choosing an Optimal Filesystem for the AI Pipeline
- Five Advantages of ONTAP AI for AI and Deep Learning
- Deep Dive into ONTAP AI Performance and Sizing

Thank You!

