https://www.nextplatform.com/2018/08/29/cascade-lake-heart-of-2019-tacc-supercomputer/



LARGE-SCALE DEEP LEARNING AS A BENCHMARK

Pradeep K Dubey

Intel Fellow and Director Parallel Computing Lab

Intel Labs



https://insidehpc.com/2018/06/univa-deploys-million-core-grid-engine-cluster-aws/



OUR EXPERIENCE SCALING DEEP LEARNING: DEEP LEARNING AT 15PF *

- Joint work between NERSC, Stanford University and Intel on Cori Phase II
- <u>Novel approach to distributed SGD</u>: synchronous within the group, asynchronous across the groups
- <u>Record scaling</u>: in terms of number of nodes collaboratively training the same model (9600 KNL)
- <u>Record peak performance</u>: ~15PF
- <u>Communication approach</u>: MLSL for intragroup communication, MPI for intergroup
- The mechanism is available in IntelCaffe/MLSL
- Real-world pattern classification problems in HEP and Climate Science



http://www.nersc.gov/users/computational-systems/cori/configuration/





WHAT SHOULD WE NEED TO CONSIDER IN A LARGE-SCALE DL BENCHMARK

- Clearly state the goal: Ranking the compute infrastructure goodness
- Scale matters 1000s of nodes
- End-to-end time-to-train to a given level of accuracy
- Consider TCO as well
- Include public cloud now offering millions of cores

- Proxy of a real-world, forward-looking application
 - Not another ImageNet 🙂
- Challenging enough DL training with spelt-out network and dataset
- Allow flexibility of data-model parallelism, inter-node communication, precision, etc.
 - All such need to be disclosed
- Open source code, reproducible results





NEURAL NETWORKS GETTING AUGMENTED





Hybrid learning of with neural networks and coupled dynamic system (PDEs) for heat dissipation and fluid dynamics [1]





Memory needed to perform Turing-complete operations. DeepMind work on differentiable memory [2]



Large embedding tables mapping sparse feature vector to dense vectors – Requires several TBs of memory [3]





FOCUS SHIFTS FROM PATTERNS TO ANOMALIES

CYBER-SECURITYREAL-TIME MONITORINGAUTONOMOUS DRIVINGDATA CENTERSSOCIAL MEDIAINTERNET OF THINGS (IOT)DATA ANALYTICSWEATHER / ASTRONOMY



Anomaly Detection Is Everywhere



NEW EVALUATION MODEL

Precision and Recall for Time Series

Nesime TatbulTae Jun LeeIntel Labs and MITMicrosofttatbul@csail.mit.edutae_jun_lee@alumni.brown.edu

Stan Zdonik Brown University sbz@cs.brown.edu

Mejbah Alam Intel Labs mejbah.alam@intel.com

Justin Gottschlich Intel Labs justin.gottschlich@intel.com

Abstract

Classical anomaly detection is principally concerned with *point-based anomalies*, those anomalies that occur at a single point in time. Yet, many real-world anomalies are *range-based*, meaning they occur over a period of time. In this paper, we present a new model that more accurately measures the correctness of anomaly detection systems for range-based anomalies, while subsuming the classical model's ability to classify point-based anomaly detection systems.

1 Introduction

Anomaly detection (AD) is the process of identifying non-conforming items, events, or behaviors. The proper identification of anomalies can be critical for many domains. Some examples are early diagnosis of illness and disease [18], threat detection for cyber-attacks [3], or safety analysis for self-driving cars [29]. Many real-world anomalies can be detected in time series data. Therefore, systems that detect anomalies should reason about them as they occur over a period of time. We call such events *range-based anomalies*, which are a subset of both contextual and collective anomalies [9]. More precisely, a *range-based anomaly* is an anomaly that occurs over a consecutive sequence of time points, where no non-anomalous data points exist between the beginning and the end of the anomaly. The standard metrics for evaluating anomaly detection algorithms today, *Recall and Precision*, have been around since the 1950s, originally formulated to evaluate document retrieval algorithms by counting the number of documents that were not [7].

Formally defined as follows, Recall and Precision are a good match for single-point AD [1] (where TP, FP, FN are the number of true positives, false positives, false negatives, respectively):

$$Recall = TP \div (TP + FN)$$
(1)

$$Precision = TP \div (TP + FP)$$
(2)

Informally, *Recall* is the rate at which a system can identify anomalies without mispredicting any anomalous events. *Precision* is the rate a system can identify anomalies without mispredicting non-anomalous events. In this sense, *Recall* and *Precision* are complementary. This characterization proves useful when they are combined, such as in the F_1 score, which is their harmonic mean. Such combinations help gauge the quality of both anomalous and non-anomalous predictions. While useful for point-based anomalies, classical recall and precision suffer from the inability to represent domain-specific time series anomalies. This has a negative side-effect on the advancement of AD systems. In particular, many time series AD systems' accuracy is being misrepresented, because point-based recall and precision are being used to measure their effectiveness for range-based anomalies. Moreover, the need to accurately identify time series anomalies is growing in importance due to the explosion of streaming and real-time systems [2, 6, 14, 23, 28, 31]. To address this, we redefine recall and precision to encompass range-based anomalies. Unlike prior work [2, 21], our mathematical

32nd Conference on Neural Information Processing Systems (NIPS 2018), Montréal, Canada.

Expressive, Flexible, Extensible

Superset of:

- Classical model
- Other state-of-the-art evaluators (NAB)
- NIPS '18 Spotlight
- Other NIPS'18 Highlight:
 - Intel for the first time in the league of >1% of accepted paper affiliations



TREND 3: MOVING BEYOND PERCEPTION Focus shifts from sparse and graph analytics

Many Big Data sets can be represented as a graph: Social networks, IP network traffic, road networks, physics models, etc. Graph analytics can reveal interesting information: detecting patterns and clusters, shortest path calculation, search problems

Differences with 'classic' HPC:

Sparse data: connection matrix has small fraction of

non-zeros

Light computations: walking through graph makes most algorithms memory bound

Data dependent: process time depends on number of neighbors, which can be highly unbalanced



