

Efficient Deep Neural Network Serving: Fast and Furious

Feng Yan^{1b}, *Member, IEEE*, Yuxiong He, *Member, IEEE*, Olatunji Ruwase, *Member, IEEE*,
and Evgenia Smirni, *Senior Member, IEEE*

Abstract—The emergence of deep neural networks (DNNs) as a state-of-the-art machine learning technique has enabled a variety of artificial intelligence applications for image recognition, speech recognition and translation, drug discovery, and machine vision. These applications are backed by large DNN models running in serving mode on a cloud computing infrastructure to process client inputs such as images, speech segments, and text segments. Given the compute-intensive nature of large DNN models, a key challenge for DNN serving systems is to minimize the request response latencies. This paper characterizes the behavior of different parallelism techniques for supporting scalable and responsive serving systems for large DNNs. We identify and model two important properties of DNN workloads: 1) homogeneous request service demand and 2) interference among requests running concurrently due to cache/memory contention. These properties motivate the design of serving deep learning systems fast (SERF), a dynamic scheduling framework that is powered by an interference-aware queueing-based analytical model. To minimize response latency for DNN serving, SERF quickly identifies and switches to the optimal parallel configuration of the serving system by using both empirical and analytical methods. Our evaluation of SERF using several well-known benchmarks demonstrates its good latency prediction accuracy, its ability to correctly identify optimal parallel configurations for each benchmark, its ability to adapt to changing load conditions, and its efficiency advantage (by at least three orders of magnitude faster) over exhaustive profiling. We also demonstrate that SERF supports other scheduling objectives and can be extended to any general machine learning serving system with the similar parallelism properties as above.

Index Terms—Deep learning, DNN serving, scheduling, parallelism, performance, analytical model, interference-aware.

I. INTRODUCTION

THE RECENT advance in Deep Neural Network (DNN) models have enabled state-of-the-art accuracy on important yet challenging artificial intelligence tasks, such as image recognition [1]–[3] and captioning [4], [5], video classification [6], [7] and captioning [8], speech recognition [9], [10],

Manuscript received May 4, 2017; revised September 22, 2017; accepted November 4, 2017. Date of publication February 21, 2018; date of current version March 9, 2018. This work is supported by NSF grant CCF-1218758, CCF-1649087, and CCF-1756013. The associate editor coordinating the review of this paper and approving it for publication was Yixin Diao. (Corresponding author: Feng Yan.)

F. Yan is with the Department of Computer Science and Engineering, University of Nevada at Reno, Reno, NV 89557 USA (e-mail: fyan@unr.edu).

Y. He and O. Ruwase are with Microsoft Research, Redmond, WA 98052 USA (e-mail: yuxhe@microsoft.com; oluwase@microsoft.com).

E. Smirni is with the Department of Computer Science, College of William and Mary, Williamsburg, VA 23187 USA (e-mail: esmirni@cs.wm.edu).

Digital Object Identifier 10.1109/TNSM.2018.2808352

and text processing [11]. These advancements by DNNs have enabled a variety of new applications, including personal digital assistants [12], real-time natural language processing and translation [13], photo search [14] and captioning [15], drug discovery [16], and self-driving cars [17].

A key driver of these recent improvements in DNN performance is the ability to train large DNN models, containing billions of neural connections, using large amounts of training data [1], [3], [16], [18]. Once trained, these big DNN models are deployed in a *serving* mode to process application inputs, such as images, voice commands, speech segments, handwritten text. However, big DNN models require significant compute cycles and memory bandwidth to process each input, and are therefore impractical to run on battery-powered and small form-factor hardware devices, such as laptops, tablets, and mobile phones. Consequently, big DNN models are typically deployed as client-server applications, with the client running on a mobile device, and the server, including the model, running as a serving system on the cloud (e.g., Cortana, Siri, and Google Now). This paper presents how to build a scalable and responsive serving systems for these large DNN models.

Like other interactive online services, such as Web search and online gaming, DNN serving requires consistently low response times to attract and retain users. Computing the answer for a user request using large DNN models may take seconds and even minutes to complete when running sequentially on a single machine.

One promising approach for reducing the DNN serving latency is to parallelize computation. There are three complementary ways to achieve large-scale parallelism in DNN serving systems. First, the DNN model, which consists of billions of neurons and connections, can be partitioned across multiple servers. Each request is processed concurrently with communications across these servers (*inter-node parallelism*). Second, at each server, a request can be further parallelized using multiple threads exploiting multicore architecture of modern hardware (*intra-node parallelism*). Finally, *multiple* requests can be processed concurrently within each multicore server to exploit service parallelism among requests (*service parallelism*). However, this service parallelism does not reduce the latency of an individual query. Nonetheless, processing multiple requests in parallel is valuable because it improves the server's throughput and potentially reduces the time that a request waits for execution.

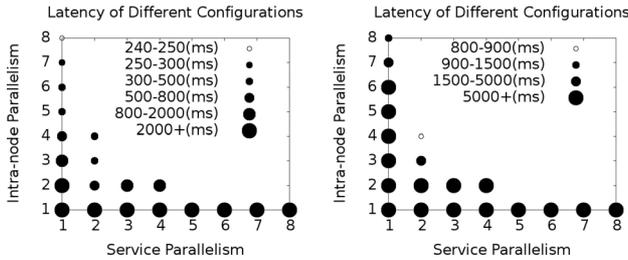


Fig. 1. Latency under low load (left plot) and high load (right plot) using different configurations (inter-node parallelism is set to 1) for ImageNet-22K.

While these parallelism techniques present opportunities to reduce DNN serving latency, deciding optimal parallelism choice is challenging. Applying parallelism degrees blindly could harm performance. For example, service parallelism may increase memory system contention to the point of prolonging request processing time; inter-node parallelism may prolong request processing if the cross-machine communication overhead exceeds the computation speedup. Figure 1 shows the latency of serving the ImageNet-22K workload [19] under different combinations of service and intra-node parallelisms on an 8-core machine (refer to Section V-A for detailed experimental setup). The left and right scatter plots represent low and high load conditions (request arrival rate) respectively. Each point represents one parallel configuration, and the size of the point indicates its latency value. The figure demonstrates that: (1) Many parallel configurations are possible, even with only 8 cores and without considering inter-node parallelism. (2) The latency difference between the best parallel configuration and the worst parallel configuration can be significant, i.e., by orders of magnitudes. This gap grows further under higher loads. (3) The latency values and the best parallel configuration changes as a function of the load.

We propose SERF, a framework for serving deep learning systems fast which integrates light-weight profiling with a queueing-based prediction model to quickly find optimal parallel configurations for DNN serving. SERF exploits three key intuitions to address the three challenges. First, it employs a dynamic scheduler that determines online the ideal parallelism configuration based on the system load. Second, as it is hard to accurately estimate the impact of a parallel configuration on the latency, SERF leverages light-weight profiling to measure workload latencies of a few key configurations on the hardware of interest. Third, instead of an exhaustive profiling over all configurations under all loads, which is unavoidably very slow and not practical, we develop the core component of SERF — a queueing-based analytical model for performance prediction — which uses only limited simple profiling (that can be done very fast) to record essential system and workload information that is used as input to the model. Using this input, the model achieves remarkably accurate predictions of the request latency of *any* parallel configuration under *any* given load, thus can be used online in a dynamic workload setting.

We implement SERF in the context of an image classification service based on the image classification module of the Adam distributed deep learning framework [3]. We stress

that SERF is not limited to the Adam architecture, but also applicable to serving systems based on other DNN frameworks (e.g., Caffe [20], Theano [21], and Torch7 [22]) as similar parallelism decisions and configuration knobs are also available there. Our current prototype includes implementations of our parallelism techniques, as well as a load generator for simulating arrival process. We evaluate its performance using a 20-machine cluster and conduct vast experiments by running several state-of-the-art classification benchmarks, including *ImageNet* [19] and *CIFAR* [23]. We demonstrate the accuracy of our queueing-based prediction model by comparing its prediction results with testbed measurements. Moreover, we show that, comparing to using static parallel configurations, SERF swiftly recommends the optimal configuration under various loads. Comparing to exhaustive profiling, SERF adapts three orders of magnitude faster under dynamic and ever-changing environments.

We also demonstrate that SERF supports different scheduling objectives, e.g., finding the minimum amount of required resources to meet a target latency SLO (service level objective) and can be extended to support any general machine learning serving system with similar characteristics as DNN serving system. We summarize the main contributions of the paper as follows: (1) We conduct a comprehensive workload characterization of a DNN serving system, highlighting the opportunities and challenges of using different parallelism techniques to reduce response latency (Section III). (2) We propose the SERF scheduling framework, which integrates lightweight profiling and queueing-based latency prediction model to find best parallel configurations effectively and efficiently (Section IV). (3) We implement SERF and evaluate it on a cluster of machines. The experimental results verify its effectiveness and efficiency (Section V).

II. BACKGROUND

DNNs consist of large numbers of neurons with multiple inputs and a single output called an activation. Neurons are connected hierarchically, layer by layer, with the activations of neurons in layer $l-1$ serving as inputs to neurons in layer l . This deep hierarchical structure enables DNNs to learn complex tasks, such as image recognition, speech recognition, and text processing.

A DNN service platform supports training and serving. DNN training is offline batch processing that uses learning algorithms, such as stochastic gradient descent (SGD) [24] and labeled training data to tune the neural network parameters for a specific task. DNN serving is instead interactive processing requiring fast response per request, e.g., within 7-10 milliseconds for speech application [25], and within 200 - 300 milliseconds even for challenging large-scale models like ImageNet-22K. It deploys the trained DNN models in serving mode to answer user requests, e.g., for a dog recognition application, a user request provides a dog image as input and receives the type of the dog as output. The response time of a request is the sum of its service time (execution time) and waiting time. An important common performance metric

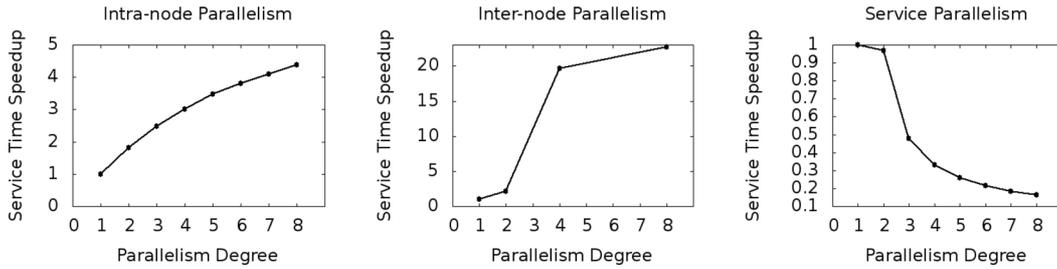


Fig. 2. Service time comparison under different parallelism techniques using ImageNet-22K. Each plot reports the speedup when increase the degree at only one parallelism (fix the other two parallelisms).

for interactive workloads is the average request response time (average latency), which we adopt in our work.

In DNN serving, each user input, which we refer to as a request, is evaluated layer by layer in a feed-forward manner where the output of a layer $l-1$ becomes the input of layer l . More specifically, define a_i as the activation (output) of neuron i in layer l . The value of a_i is computed as a function of its J inputs from neurons in the preceding layer $l-1$ as: $a_i = f((\sum_{j=1}^J w_{ij} \times a_j) + b_i)$, where w_{ij} is the weight associated with the connection between neuron i at layer l and neuron j at layer $l-1$, and b_i is the bias term associated with neuron i . The activation function f , associated with all neurons in the network, is a pre-defined non-linear function, typically a sigmoid or hyperbolic tangent. Therefore, for a given request, its main computation at each layer l is a matrix-vector multiplication of the weight of the layer with the activation vector from layer $l-1$ (or the input vector if $l=0$).

Inter-node, intra-node, and service-level parallelisms are well-supported among various DNN models and applications [1], [3], [26]. Inter-node parallelism partitions the neural network across multiple node/machines, with activations of neural connections that cross node/machine boundaries being exchanged as network messages. Intra-node parallelism uses multi-threading to parallelize the feed-forward evaluation of each input image using multiple cores. As the computation at each DNN layer is simply a matrix-vector multiplication, it can be easily parallelized using parallel libraries such as OpenMP [27] or TBB [28] by employing a parallel for loop. Service-level parallelism is essentially admission control that limits the maximum number of concurrent running requests. We define a *parallelism configuration* as a combination of the intra-node parallelism degree, inter-node parallelism degree, and maximum allowed service parallelism degree. Note the service parallelism is defined as a maximum value instead of the exact value due to the random request arrival process, e.g., at certain moments, the system may have less requests than the defined service parallelism degree.

III. WORKLOAD CHARACTERIZATION

In this section, we present comprehensive workload characterization that shows the opportunities and challenges of using the various parallelism techniques to reduce DNN serving latency, as well as their implications on the design of SERF. We make four key observations: (1) Parallelism impacts service time in complex ways, making it difficult to model service

times without workload profiling. (2) DNN workloads have homogeneous requests, i.e., service times under the same parallelism degree exhibit little variance, which allows SERF to measure request service time with affordable profiling cost. (3) DNN workloads exhibit interference among concurrent running requests, which motivates a new model and solution of SERF. (4) DNN workloads show load-dependent behavior, which indicates the importance of accurate latency estimation and parallel configuration adaptation according to system load.

We present workload characterization results of two well-known image classification benchmarks, CIFAR-10 [2] and ImageNet-22K [19], on servers using Intel Xeon E5-2450 processors. Each processor has 8 cores, with private 32KB L1 and 256KB L2 cache, and shared 20MB L3 cache. The detailed experimental set up for both workloads and hardware is provided in Section V.

A. Impact of Parallelism on Service Time

Modeling the impact of parallelism on DNN serving without workload profiling is challenging because parallelism has complex effects on the computation and communication components of request service time (shown in Figures 2 and 3).

Figure 2 shows the DNN request service speedup for different degrees of intra-node, inter-node, and service parallelism. For intra-node parallelism, the speedup is close to linear up to 3 cores, but slows down beyond 4 cores. This effect is due to the limited memory bandwidth. When the total memory bandwidth demands are close to or exceed the available bandwidth, the bandwidth per core reduces, decreasing speedup. For inter-node parallelism, increasing the parallelism degree from 1 to 2 yields a 2X service time speedup because the computation time, which is dominant, is halved, while communication time grows marginally; increasing from 2 to 4 results in super-linear speedup due to caching effects, as the working set fits in the L3 cache; increasing from 4 to 8 results in smaller speedup increase as communication starts to dominate service time. For service parallelism, parallelism degrees > 2 result in increased service time due to memory interference among concurrently serviced requests. These results are indicative of the impact of different parallelism on service time. Speedups can vary a lot, depending on many factors, including DNN size, the ratio of computation and communication, cache size, memory bandwidth.

Figure 3 demonstrates the relationship between inter-node and intra-node parallelism: the results indicate that the degree

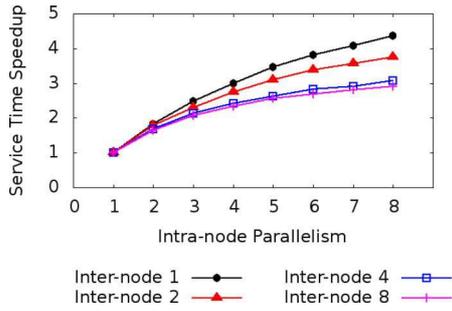


Fig. 3. Relationship between inter-node and intra-node parallelism using ImageNet-22K.

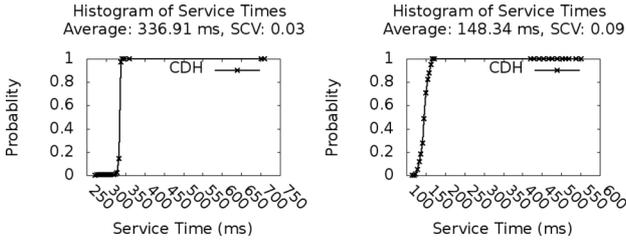


Fig. 4. CDH (Cumulative Data Histogram) of service times. The left plot is with parallelism degree tuple (2, 1, 4) and the right plot is with (4, 4, 2).

of one parallelism technique can affect the behavior of another. More precisely, intra-node parallelism speedup depends on the degree of inter-node parallelism: speedup reduces with larger inter-node parallelism. This is because communication time is increasingly the dominant portion of service time with larger degrees of inter-node parallelism, therefore the computation time improvements of intra-node parallelism become less important to overall service time.

In summary, since parallelism efficiency depends on various factors (e.g., workload and hardware properties) and since one parallelism technique can affect the behavior of others, it is difficult to accurately model service time. SERF circumvents this by incorporating workload profiling to predict request service time.

B. Homogeneous Requests

We observe that for a given *parallelism degree tuple*,¹ defined as (service parallelism degree, inter-node parallelism degree, intra-node parallelism degree), the service times of DNN requests exhibit very little variance because the same amount of computation and communication is performed for each request. Thus, we refer to DNN requests as being homogeneous. Figure 4 shows two examples corresponding to two representative cases of parallelism degrees. The first example as shown in the left plot of Figure 4 is with parallelism degree tuple of (2, 1, 4), where the majority of requests are in the range of 330ms to 340ms and the SCV (squared coefficient of variation) is only 0.03. The second example as shown in the right plot of Figure 4 is under parallelism (4, 4, 2),

¹Note that parallelism degree tuple is different from parallelism configuration. In parallelism degree tuple, each parallelism is set exactly to the degree value while in parallelism configuration, max service parallelism is an admission policy that defines the maximum allowed degree of service parallelism.

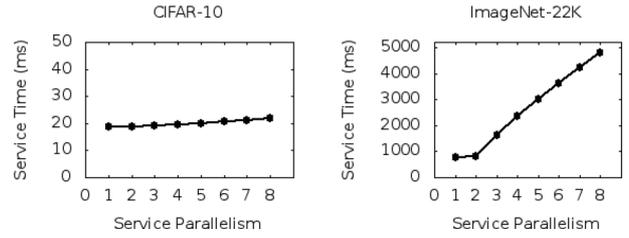


Fig. 5. Service time comparison with different number of concurrent requests.

where most requests are in the range of 130ms to 160ms with the SCV of 0.09. The slightly larger variance can be attributed to variations in the cross-machine communication delays caused by inter-node parallelism. The magnitude of these variations is consistent with what is normally expected in computer communication systems while running a request multiple times [29].

This unique property of homogeneous requests for DNN workloads empowers lightweight profiling: the cost of measuring the service time is low, i.e., for a given parallelism degree tuple, running one or a few input requests is sufficient. In comparison, many other online services have requests with heterogeneous demands [30], [31] and require to execute many more input samples to collect service time distributions.

C. Interference Among Concurrent Requests

For small DNNs like CIFAR-10 (the left plot of Figure 5), request service time remains almost constant when running requests concurrently under different service parallelism degrees, because there is little interference among requests due to cache/memory contention. The interference becomes more obvious for large DNNs. The right plot in Figure 5 shows the request service time of ImageNet-22K when running different number of requests. It is clear that when running more than 2 requests concurrently, the interference becomes severe. To explain performance interference, it is important to understand the working set of DNN serving that comprises activations and weights of the neural connections (the core operation is a matrix-vector multiplication of the weight matrix and the activation vector). Activations are derived from request input, while weights represent the model parameters and are shared by all requests. When there are no more than 2 concurrent requests, the working sets of both fit into L3 cache. If more than three requests run concurrently, then the footprint of activations increases and the aggregate working set no longer fits in the L3 cache, resulting in more L3 cache misses, thus prolonging the request service time. This is also why large DNNs like ImageNet-22K are more likely to have interference than small ones, such as CIFAR-10. Note that even with the compression techniques and more powerful hardware, the working set cannot always fit into L3 cache as there can be more complex and larger models/data being used, so such interference behavior is ubiquitous.

Interference makes modeling average service time and waiting time for a given parallelism configuration much more challenging. In particular, under the same parallelism configuration, the number of running requests can vary from 0 to the

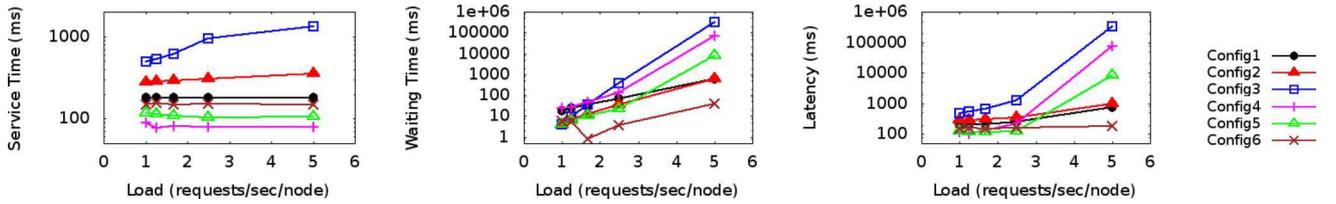


Fig. 6. Service time, waiting time, and latency under different loads using different configurations for ImageNet-22K.

TABLE I
PARALLEL CONFIGURATIONS

Config.	Service	Inter-node	Intra-node
Config1	1	1	8
Config2	2	1	4
Config3	4	1	2
Config4	1	4	8
Config5	2	4	4
Config6	4	4	2

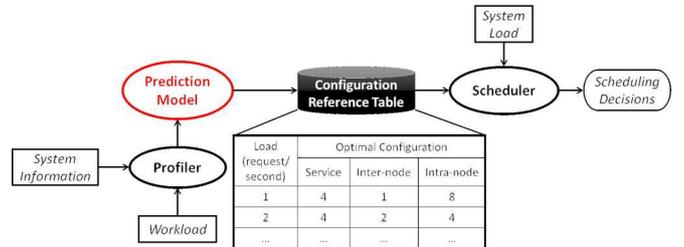


Fig. 7. Overview of SERF.

maximum service parallelism of the configuration. Therefore, the service time of a particular request depends on the number of concurrent running requests at the moment of its execution, and the average service time depends on the probability distribution of the concurrency levels. The waiting time estimation is even more complex. The existing queuing and scheduling models [32] are no longer applicable as they assume independence among requests: request service time remains constant regardless of the number of concurrent requests. This property of DNN motivates us to develop new model and solution of SERF to accurately model the waiting time and latency impact of interference.

D. Load-Dependent Behavior

In serving systems, load (request arrival rate) changes dynamically over time. For a given parallel configuration, both request service time and waiting time could change under different loads. To illustrate the load-dependent behavior of different parallelism approaches, we use 6 distinctive configurations and conduct experiments under different load levels, see Table I.

The left plot in Figure 6 shows the service time of using the 6 configurations under different loads, the middle plot in Figure 6 shows their waiting time, and the right plot in Figure 6 shows their latency. The results demonstrate that for the same configuration, service time, waiting time, and latency can vary under different loads. Therefore, the ability to estimate the latency impact according to the load and a scheduler that can change the parallel configurations based on load are two necessary and important features.

IV. SERF: A FRAMEWORK FOR DNN SERVING

In this section, we present the scheduling framework SERF. SERF applies a hybrid approach that integrates lightweight profiling and a queueing-based prediction model to find best parallel configurations for any given load (request arrival rate) effectively and efficiently, achieving the benefits of both empirical and analytical methods. We first discuss the scheduling

objective and give an overview of SERF (Section IV-A). Then we answer the two important questions raised in the Introduction: (1) What should be profiled for accuracy yet can be profiled quickly (Section IV-B)? (2) How to model the rest and predict request latency (Section IV-C)? Finally, we discuss how to use the prediction results to change the parallelism configurations online with varying loads (Section IV-D) and its support for other scheduling objectives (Section IV-E).

A. Overview

Scheduling Objective: Common objectives for scheduling interactive serving systems are (1) to minimize response latency using a given amount of resources [30], [31] or (2) to minimize resource consumption while meeting latency SLO [33], [34]. Our scheduling framework supports both. We use the first objective of minimizing response latency as example to develop the key components of the framework, then we discuss how to extend the proposed approach for the second objective. We choose to optimize average latency because DNN requests are homogeneous and have similar service time, reducing average latency also reduces the tail latency.

Framework Overview: Figure 7 presents an overview of SERF, which consists of three main modules: prediction model, profiler, and scheduler. The modules are connected by the configuration reference table, which maps different load levels (represented by request arrival rate) to their corresponding best parallel configurations. For example, at arrival rate of 2 requests/second, the best configuration is with a max service parallelism 4, inter-node parallelism of 2, and intra-node parallelism of 4. The profiler takes the system information (e.g., the number of machines and cores, and workload) as input and conducts lightweight profiling and feeds the profiling results to the prediction model. The prediction model is the key component of the framework. It utilizes the profiling results to predict the latency of all combinations of parallelism under different load levels and populates the configuration reference table. This table only needs to be built once, provided that

DNN workload characteristics and system hardware remain the same. The scheduler uses the current system load as index to search the configuration reference table, find and adapt to the best parallel configurations.

B. Profiler

An easy but inefficient way to achieve the scheduling objective is via exhaustive profiling: execute all possible parallelism configurations for all possible loads and find the best parallel configuration for each load. The shortcoming of such exhaustive profiling is its high cost. Assuming that there are P different configurations and there are L load levels, one needs to conduct $P \times L$ profiling experiments. In addition, measuring average latency requires a relatively long time span (to measure enough samples) to achieve statistical stability due to the stochastic queuing behaviors. Experimenting with lighter load levels requires even longer time for profiling because the large idle intervals between requests increase the duration of the experiment. Let T be the average cost to achieve statistical stability in profiling, which makes the overall cost of exhaustive profiling $P \times L \times T$.

SERF conducts lightweight profiling by measuring the *request service time* for each parallelism degree tuple of (service parallelism degree, inter-node parallelism degree, intra-node parallelism degree). For example, with the tuple (2, 4, 3), we measure the request service time by running two requests concurrently, each request across 4 server nodes and with 3 cores on each server node. Let E denote the cost of profiling request service time for a given parallelism degree tuple, the total profiling cost of SERF is $P \times E$, where P is the total number of parallelism degree combinations. The profiling of SERF has two key differences compared to exhaustive profiling, resulting in significantly lower profiling cost: (1) SERF measures the request *service time* instead of latency, and (2) SERF measures each parallelism degree tuple instead of each parallel configuration. Benefit of these profiling choices is two-fold: (1) the service time of different parallelism configurations under different loads can be computed by SERF, saving a multiplicative cost factor along the load dimension L . (2) As requests have almost deterministic service time under the same parallelism degree tuple and profiling the service time is independent of the queuing delays, a few profiling samples are sufficient, i.e., the value of E is small. In contrast, exhaustive profiling measures latency for each parallelism configuration, which requires running many samples to achieve statistical stability for queuing delays, i.e., T is much more costly than E , by up to 3 orders of magnitude. Therefore, SERF profiling is much more efficient than exhaustive profiling, and $P \times E \ll P \times L \times T$. We feed these profiling results to the prediction model of SERF to estimate the latency under different load levels, which is introduced next.

C. Queuing-Based Prediction Model

We develop a queuing model that takes profiling results as input and predicts request latency under different load and parallelism configurations. The key challenge and novelty of the model is its interference-awareness, effectively quantifying

the latency impact of request interference due to cache and memory contention.

1) *Problem Formulation*: We define the problem as predicting DNN request latency for any given parallel configuration under any given load. We denote parallelism configuration with (*maximum service parallelism* $C_{service}$, *inter-node parallelism* C_{inter} , and *intra-node parallelism* C_{intra}). The inputs of the model are:

- Load in terms of inter-arrival rate: λ , here we assume Poisson arrivals for a *short period*, i.e., exponential inter-arrival times with mean rate λ , which is typical for online services [35], [36]. Such assumption does not contradict the bursty and long-range dependence characteristics in [37]. SERF continuously monitors the incoming workload and periodically updates its observed load (arrival rate).
- Profiling results: μ_i ($i = 1 \dots c$) represents the average service rate when i requests are running concurrently, i.e., the average service rate of the parallelism degree tuple (i, C_{inter}, C_{intra}).

The output of the model is the average latency for the parallelism configuration under any given load.

We model DNN serving as an interference-aware deterministic service process and formulate the problem as a $M/D_{interf}/c$ queue. Here, M represents exponential inter-arrival times. D_{interf} represents two distinctive properties of DNN workload: (1) **D**eterministic service times, modeling homogeneous requests that exhibit little service time variance for any given parallelism degree tuple (as shown in Section III-B). (2) **I**nterference-awareness, modeling the interference among requests due to cache and memory contention (as shown in Section III-C). c stands for the maximum service parallelism, equal to $C_{service}$.

2) *Technical Challenges and Key Ideas*: The $M/D_{interf}/c$ queue does not have a closed-form solution. In fact, even for the simpler problems: the interference-oblivious $M/D/c$ queue that assumes deterministic service time without any interference among concurrent running requests, or interference-aware $M/M_{interf}/c$ queue that assume exponential distributed service times with interference among concurrent running requests, there is no closed-form solution. Intuitively, one may want to use $M/M_{interf}/c$ queue, $M/D/c$ queue, or $M/M/c$ queue to approximate the $M/D_{interf}/c$ queue, but such approximation has the potential of achieving bad accuracy. To illustrate why these simpler approaches can not model DNN workload, we implement these approximation methods and conduct experiments using ImageNet-22K. We compare the latency results of best configurations under different loads between testbed measurements and prediction results from $M/M_{interf}/c$ queue, $M/D/c$ queue, and $M/M/c$ queue in Figure 8. The results clearly shows that the prediction from these approaches is poor. This large discrepancy shows the importance of incorporating interference and deterministic service times into SERF prediction model and solution.

Our solution is inspired by Cosmetatos' approximation [38] that estimates $M/D/c$ model using the $M/M/c$ model with adjustment and correction, where $M/M/c$ model is a standard multi-server queue model with Poisson arrival and exponential

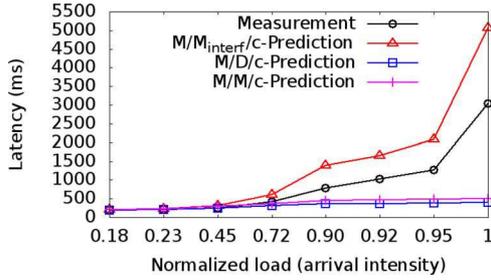


Fig. 8. Latency comparison of best configurations between measurement and standard prediction under different load.

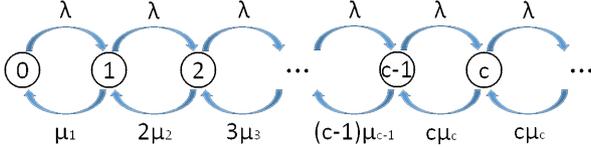


Fig. 9. State transition diagram for the $M/M_{interf}/c$ queue. Each state represents the number of requests in the node.

service time. We extend the approximation approach to the interference-aware case and solve $M/D_{interf}/c$ queue in two steps. (1) Solve $M/M_{interf}/c$ queue that has interference-aware exponential service time. (2) Utilize the approximation method proposed in Cosmetatos' approximation to adjust the results of $M/M_{interf}/c$ queue to approximate the $M/D_{interf}/c$ queue. We estimate the waiting time and service time separately. Latency is estimated as the sum of these two measures.

3) *Solving $M/D_{interf}/c$ Queue: Waiting time estimation:* We follow the two steps described in Section IV-C2 to solve for the waiting time.

(1) *Solving $M/M_{interf}/c$ queue:* Recall that μ_i ($i = 1 \dots c$) is provided by profiling and represents the average service rate when i requests are concurrently running, p_i is the probability of i requests in the system. Let $\rho_i = \frac{\lambda}{\mu_i}$ and $\rho = \frac{\rho_c}{c} = \frac{\lambda}{c \cdot \mu_c}$. Based on the state transition diagram shown in Figure 9 and global balance equations [39], we obtain:

$$p_n = \begin{cases} \frac{\prod_{i=1}^n \rho_i}{n!} \cdot p_0 & (0 \leq n \leq c-1) \\ \frac{\rho^{n-c} \cdot \prod_{i=1}^c \rho_i}{c!} \cdot p_0 & (n \geq c), \end{cases} \quad (1)$$

where p_n is the steady-state probability of state n , which represents n requests in the system (sum of requests in the queue and in service). p_0 represents the probability that the system is idle, i.e., no request is in the system. Since all probabilities sum to 1 and $\sum_{k=0}^{\infty} \rho^k = \frac{1}{1-\rho}$ for $\rho < 1$:

$$\begin{aligned} \sum_{k=0}^{\infty} p_k &= p_0 \cdot \left(1 + \sum_{k=1}^{c-1} \frac{\prod_{i=1}^k \rho_i}{k!} + \frac{\prod_{i=1}^c \rho_i}{c!} \cdot \sum_{k=c}^{\infty} \rho^{k-c} \right) \\ &= p_0 \cdot \left(1 + \sum_{k=1}^{c-1} \frac{\prod_{i=1}^k \rho_i}{k!} + \frac{\prod_{i=1}^c \rho_i}{c! \cdot (1-\rho)} \right) = 1, \quad (2) \end{aligned}$$

where $\rho < 1$.

Let $H = 1 + \sum_{k=1}^{c-1} \frac{\prod_{i=1}^k \rho_i}{k!} + \frac{\prod_{i=1}^c \rho_i}{c! \cdot (1-\rho)}$, then:

$$p_0 = H^{-1}. \quad (3)$$

Assume that $L_q(\lambda)$ is the average number of requests waiting in the queue, by definition we have:

$$L_q(\lambda) = \sum_{k=c}^{\infty} (k-c) \cdot p_k, \quad (4)$$

together with Eq. (1) and based on $\sum_{k=0}^{\infty} k \rho^k = \frac{\rho}{(1-\rho)^2}$ for $\rho < 1$, we have:

$$\begin{aligned} L_q(\lambda) &= \frac{p_0 \cdot \prod_{i=1}^c \rho_i}{c!} \cdot \sum_{k=c}^{\infty} (k-c) \cdot \rho^{k-c} \\ &= \frac{p_0 \cdot \prod_{i=1}^c \rho_i}{c!} \cdot \frac{\rho}{(1-\rho)^2}, \end{aligned} \quad (5)$$

where $\rho < 1$. Using Little's law [40], the waiting time in the queue can be computed as:

$$W_q^{M/M_{interf}/c}(\lambda) = \frac{L_q(\lambda)}{\lambda} = \frac{p_0 \cdot \prod_{i=1}^c \rho_i}{\lambda \cdot c!} \cdot \frac{\rho}{(1-\rho)^2}, \quad (6)$$

(2) *Approximating $M/D_{interf}/c$ using $M/M_{interf}/c$:* Cosmetatos' approximation proposed in [38] states that the waiting time in the queue can be approximated as:

$$W_q^{M/D/c} \approx \frac{1}{2} (1 + f(c) \cdot g(\rho)) \cdot W_q^{M/M/c}, \quad (7)$$

where

$$f(c) = \frac{(c-1) \cdot (\sqrt{4+5c} - 2)}{16c}, \quad (8)$$

$$g(\rho) = \frac{1-\rho}{\rho}, \quad (9)$$

$\rho = \frac{\lambda}{c \cdot \mu}$, λ is the average arrival rate, and μ is the average service rate. This approximation can be adjusted for the interference-aware case: we use the $M/M_{interf}/c$ queue with the same correction terms as $f(c)$ and $g(\rho)$ as above to approximate the $M/D_{interf}/c$ queue as follows (using Eq. (6) and Eq. (7)):

$$W_q^{M/D_{interf}/c}(\lambda) \approx \frac{1}{2} (1 + f(c) \cdot g(\rho)) \cdot \frac{p_0 \cdot \prod_{i=1}^c \rho_i}{\lambda \cdot c!} \cdot \frac{\rho}{(1-\rho)^2}. \quad (10)$$

Service time estimation: Although service time under the same parallelism degree tuple is deterministic and can be profiled, the service time under a given parallel configuration could change with load and needs to be predicted. This is because of the random requests arrival process and interference, e.g., at different moments, the system may have different number of concurrent running requests (ranging from 0 to the defined maximum service parallelism), which

results in different interference and therefore different service times. We use the PASTA (Poisson Arrivals See Time Averages) property [41] to compute the average service time $S^{M/D_{interf}/c}(\lambda)$ under arrival rate λ as follows:

$$\begin{aligned} S^{M/D_{interf}/c}(\lambda) &= \frac{1}{\mu_1} \cdot p_0 + \frac{1}{\mu_2} \cdot p_1 + \frac{1}{\mu_3} \cdot p_2 + \dots \\ &\quad + \frac{1}{\mu_c} \cdot p_{c-1} + \frac{1}{\mu_c} \cdot \sum_{i=c}^{\infty} p_i \\ &= \sum_{i=1}^c \frac{p_{i-1}}{\mu_i} + \frac{p_0 \cdot \prod_{i=1}^c \rho_i}{\mu_c \cdot c! \cdot (1 - \rho)}. \end{aligned} \quad (11)$$

Latency estimation: The average latency $W^{M/D_{interf}/c}$ equals to the average time spent in waiting in queue $W_q^{M/D_{interf}/c}$ plus the average time spent in execution $S^{M/D_{interf}/c}$:

$$\begin{aligned} W^{M/D_{interf}/c}(\lambda) &\approx \sum_{i=1}^c \frac{p_{i-1}}{\mu_i} + \frac{p_0 \cdot \prod_{i=1}^c \rho_i}{\mu_c \cdot c! \cdot (1 - \rho)} \\ &\quad + \frac{1}{2} (1 + f(c) \cdot g(\rho)) \cdot \frac{p_0 \cdot \prod_{i=1}^c \rho_i}{\lambda \cdot c!} \\ &\quad \times \frac{\rho}{(1 - \rho)^2}. \end{aligned} \quad (12)$$

In the above formula, recall that μ_i is an input from profiling, λ is affected by inter-node parallelism C_{inter} as it defines how many machines to serve each request, c equals to the maximum allowed service parallelism $C_{service}$, and the intra-node parallelism is restricted by F/c , where F is the number of cores in a node. Therefore, for a given system and workload, latency can be computed under different combinations of service parallelism, inter-node parallelism, and intra-node parallelism. Eq. (12) is used to populate the configuration reference table that is the core of SERF.

The above solution is derived for a single serving unit (with C_{inter} number of machines). For a cluster, the cluster can be divided into serving units based on the inter-node parallelism C_{inter} , e.g., for a cluster with N machines, there are N/C_{inter} units and each unit has an arrival rate of $\lambda = \frac{\lambda_{all}}{N/C_{inter}}$, where λ_{all} is the request arrival rate at the cluster.

D. Scheduler

The scheduler takes the current system load as input, searches the configuration reference table, finds and adapts to the best parallelism configuration. To enable quick configuration switching, the entire DNN model is pre-installed on each server, and each input is sent to the server with a mapping of servers to input partitions. This informs the server of which partition of the DNN model to use to process the input, and which servers to communicate with for cross-machine neural connections.

To sum up, we explore two distinctive properties of DNN workload — homogeneous requests with interference — to develop SERF. SERF combines lightweight profiling with an interference-aware queuing model to predict DNN serving

latency. It finds the best parallel configuration for any given load and then deploys a dynamic scheduler to adapt to varying loads online nearly instantly.

E. Meeting SLO With Minimum Resources

This section demonstrates another usage scenario of SERF, meeting SLO (service level objective) with minimum resources, which is an important objective for many online services and cloud applications [33], [34]. To meet latency SLOs for a given load (request arrival rate), SERF uses the prediction model to compute the minimum resources (e.g., the number of machines) required by each parallelism configuration and chooses the best configuration that requires the least resources. To compute the minimum resource for a given load, SERF searches through all possible parallelism configurations starting from the least amount of total resources. If none of the parallelism configurations can achieve latency SLO, we increase the amount of total resources and search again all possible parallelism configurations. We repeat such process until we find a configuration that can achieve latency SLO. Similar process is repeated for other load levels for finding the optimal configurations. Then we build the configuration reference table with these optional configurations so that SERF can adapt the parallelism configuration with the change of load.

V. EXPERIMENTAL EVALUATION

We present experimental results demonstrating how SERF improves DNN serving performance with respect to minimizing response latency. Specifically, we evaluate the following properties of SERF: (i) identify configurations under the same load level, (ii) latency change trend when load level changes, (iii) accuracy of the latency prediction model, (iv) adaptability to load dynamism compared to a static configuration, (v) efficient best configuration search compared to exhaustive profiling, (vi) another usage scenario of SERF on minimizing resources for meeting latency SLO, and (vii) how SERF works in other machine learning serving systems.

A. Experimental Setup

System Overview: We prototyped SERF based on the Adam distributed DNN system [3], which supports service parallelism through admission control, intra-node parallelism using OpenMP [27], and inter-node parallelism by partitioning the model across different machines. In order to quickly switch the configurations, the entire model parameter is pre-installed on each server, and each input is augmented to the server with a mapping of servers to input partitions. As most distributed DNN serving platforms support part or all of these parallelisms, SERF can be used in other systems as well.

Workload: We evaluate SERF using 3 popular image recognition tasks of varying complexity with Poisson request arrivals:

- CIFAR-10 [2]: classifies 32×32 color images into 10 categories. The DNN is moderately-sized, containing about

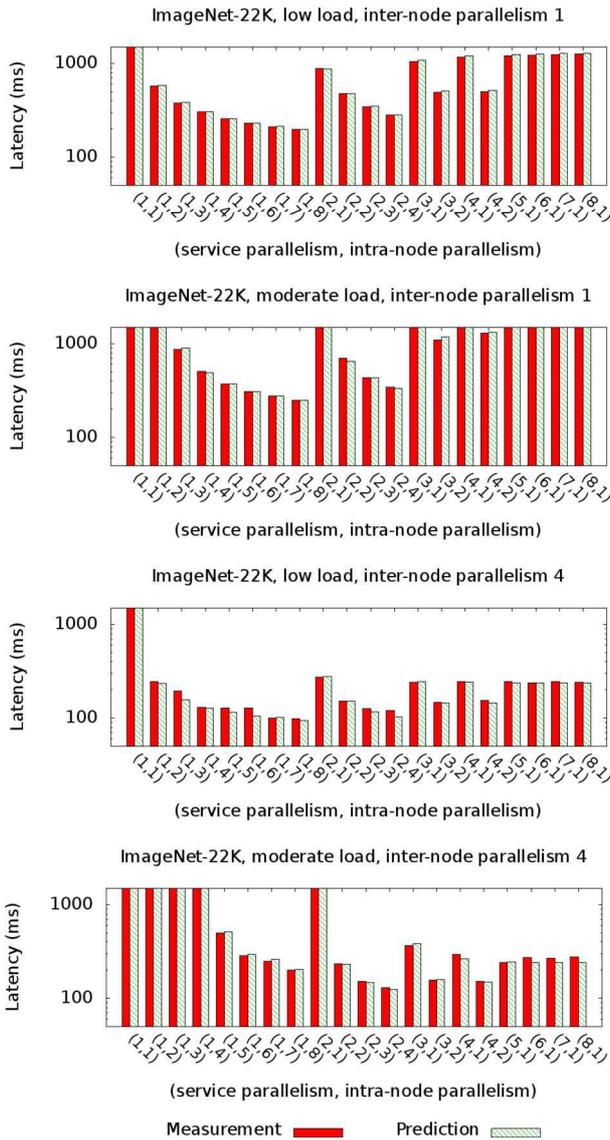


Fig. 10. Identify configurations under the same load level for different service and intra-node parallelism of ImageNet-22K. Inter-node parallelisms of 1 and 4 are demonstrated (top two, bottom two plots).

28.5 million connections in 5 layers: 2 convolutional layers with pooling, 2 fully connected layers, and a 10-way output layer.

- ImageNet-1K [19]: classifies 256×256 color images into 1,000 categories. The DNN is moderately large, containing about 60 million connections in 8 layers: 5 convolutional layers with pooling, 3 fully connected layers, and a 1,000-way output layer [2].
- ImageNet-22K [19]: the largest ImageNet task, which is to classify 256×256 color images into 22,000 categories. This DNN is extremely large, containing over 2 billion connections in 8 layers: 5 convolutional layers with pooling, 3 fully connected layers, and a 22,000-way output layer [3].

Hardware Environment: Experiments are run on a computing cluster of 20 identically configured commodity machines, communicating over Ethernet through a single 10Gbps

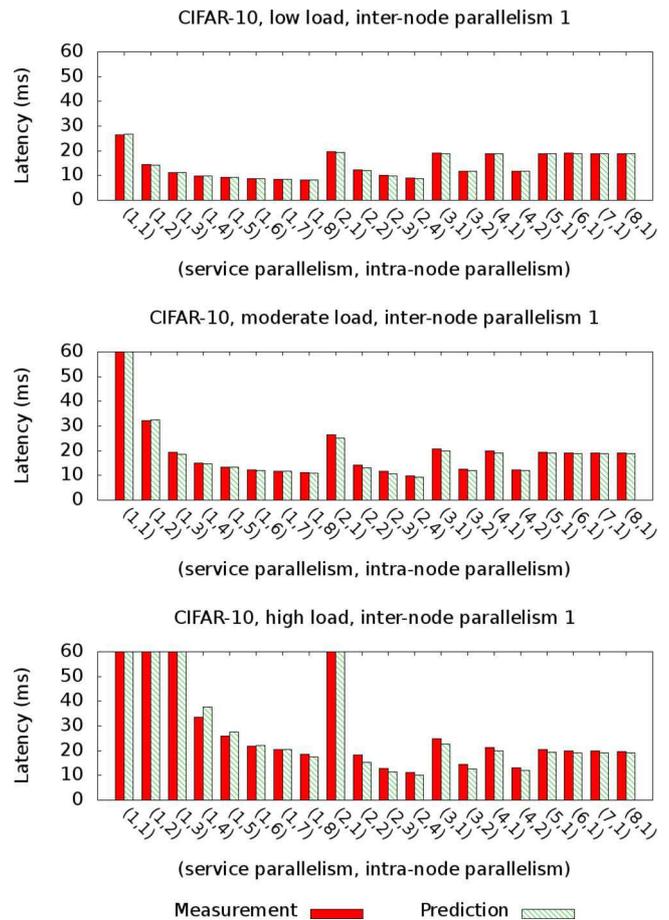


Fig. 11. Identify configurations under the same load level for different service and intra-node parallelism of CIFAR-10. Inter-node parallelisms of 1 under three different loads are demonstrated: low (top), moderate (middle), high (bottom).

(bidirectional) NIC. Each machine is dual-socket, with an Intel Xeon E5-2450 processor of 8 cores running at 2.1GHz on each socket. Each machine has 64 GB of memory and a 268.8 GFLOP/s SIMD FPU.

B. Identify Configurations Under the Same Load Level

We visualize the configuration identification by showing the results of all parallel configurations under the same load level. We select load levels of 1 and 2.5 requests/second/node as demonstration cases, representing low and moderate load levels respectively. Results are presented in Figure 10. We organize the results based on the inter-node parallelism of 1 (upper two plots) and 4 (bottom two plots). For each inter-node parallelism, we show the two selected load cases (low and moderate). For each load case, we enumerate all combinations of service parallelism and intra-node parallelism. Note that combinations are restricted by the total number of cores in each server node, i.e., 8 cores in our testbed. Figure 10 indicates that the prediction is consistently accurate across all parallelism configurations.²

²We show the latency up to 1500ms for a clear presentation. A configuration with latency beyond 1500ms is not usable in practice.

TABLE II
COST COMPARISON BETWEEN EXHAUSTIVE PROFILING AND SERF FOR DIFFERENT BENCHMARKS

Benchmark	Method	# of configs	# of load levels	# of profile exp to run	Each profile time (min)	Total time (min)
ImageNet-22K	Exhaustive	80	10	800	34.50	27600
	SERF	80	0	80	0.07	5.52
ImageNet-1K	Exhaustive	40	10	400	10.83	4332
	SERF	40	0	40	0.02	0.87
CIFAR-10	Exhaustive	20	10	200	0.36	72
	SERF	20	0	20	7.19×10^{-4}	1.44×10^{-2}

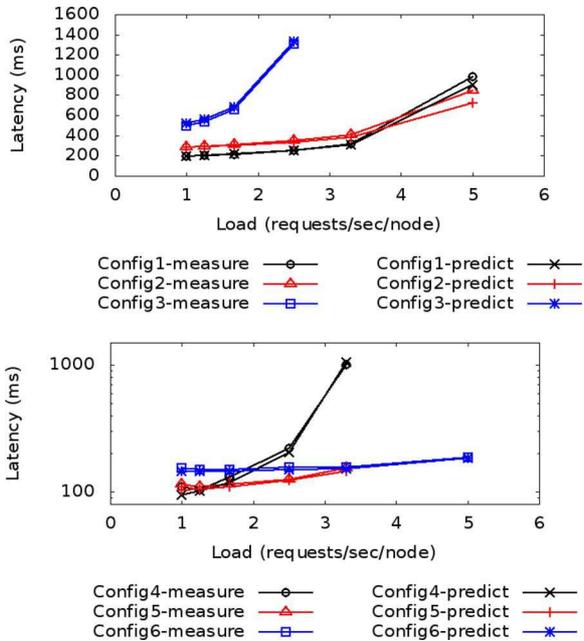


Fig. 12. Latency trend change prediction under different loads for ImageNet-22K. Inter-node parallelism of 1 and 4 are considered (top, bottom plots).

We also show the results of a much smaller application CIFAR-10, see Figure 11. Because CIFAR-10 is very small, we set inter-node parallelism to 1 and show three different load levels (low (25 requests/second/node), moderate (67 requests/second/node), and high (100 requests/second/node)). Due to the interest of space, we skip the results for ImageNet-1K, which is a similar but smaller workload compared to ImageNet-22K, but we emphasize that we can correctly identify configurations across all cases.

C. Latency Change Trend When Load Level Changes

We show here the prediction results of our model against the testbed measurement results for the parallel configurations defined in Table I. Figure 12 shows that model and experimental data are in excellent agreement. Even though the average latency of different parallel configurations have quite different trends, the prediction results can capture well these trends. Note that some parallel configurations result in system overload, i.e., an unstable system where the queue length accumulation grows to infinity and system utilization is 100%. No model can capture the behavior of an unstable system, so we do not plot points that approach the asymptote, e.g., config-3 in Figure 12.

D. Accuracy of Latency Prediction Model

This section evaluates the accuracy of the latency prediction model, based on Eq. (12), by comparing predicted values to measured values. Figure 13 shows for each workload the average and distribution of prediction errors for all relevant prediction cases. A relevant prediction case is a combination of a parallel configuration that has performance impact for a workload and a load level. For example, CIFAR-10 has 20 parallel configurations because inter-node parallelism degrees > 1 do not make sense for its small size. The larger ImageNet-1K and ImageNet-22K have 40 and 80 parallel configurations because inter-node parallelism degrees of up to 2 and 4 are relevant, respectively. For each benchmark we consider 10 load levels evenly spread across low load to high load, so that there are 200, 400, and 800 relevant prediction cases for CIFAR-10, ImageNet-1K, and ImageNet-22K, respectively. The results show that the prediction is accurate and the errors are insignificant: the average error is 2-4%, the 90th percentile is $< 10\%$, and the 95th percentile is $< 12\%$.

E. Benefits Over Exhaustive Profiling

We evaluate SERF here against exhaustive profiling for identifying the best parallel configurations under different load levels. The experimental results verified both SERF and exhaustive profiling *always* correctly identifies the best configuration. However, the cost of SERF is significantly lower than exhaustive profiling. Assume that the system has 80 different parallel configurations and the performance reference table has 10 entries (e.g., 10 different load levels). SERF requires only 80 quick profiling experiments while exhaustive profiling requires 800 expensive profiling experiments to build the performance reference table. The time for each profiling experiment and the total time to build the performance reference table is shown in Table II. Note SERF requires much less time for each profiling experiment because it only samples the service time and the service time is deterministic without load impact (i.e., sample the service time of only 10 requests) while each exhaustive profiling experiment needs to measure the average latency, which needs many samples to achieve statistical stability (e.g., when measuring latency less than 5000 sample requests, the results become very unstable). The results suggest that the time cost of SERF is more than 3 orders of magnitudes lower than exhaustive profiling, and the time savings grows with the size of the DNN workload and the number of performance reference table entries. Even if compared to lightweight profiling, e.g., only do profiling under high load, the cost of SERF is still more than 2 orders of magnitudes lower.

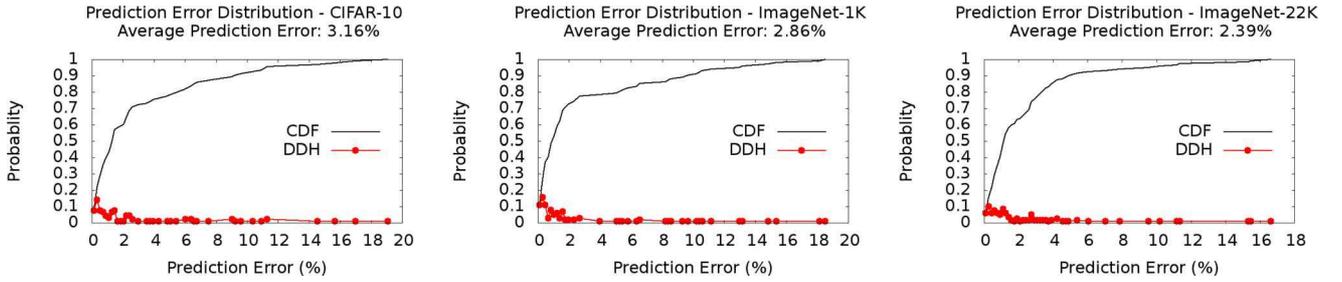


Fig. 13. CDF (Cumulative Distribution Function) and DDH (Data Density Histogram) of prediction errors for different workloads.

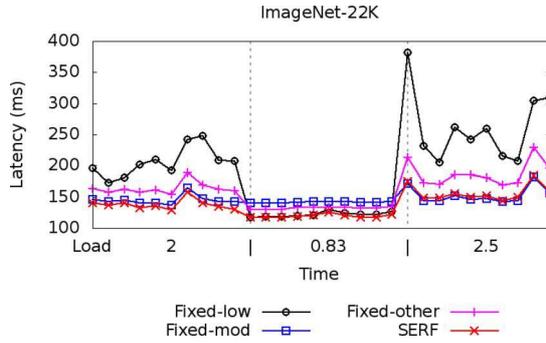


Fig. 14. Latency comparison under dynamic load environment for different scheduling approaches.

F. Benefits Over Static Configuration

Request arrival rate and system load changes dynamically for online services [36]. In this section, we demonstrate how SERF outperforms fixed configurations by adapting to load changes. We use three baseline cases for comparison. Fixed-low is a best configuration in low load and Fixed-mod is a best configuration in moderate load. Fixed-other is another configuration that performs better than Fixed-low in moderate loads and better than Fixed-mod in low loads. We compare the performance of SERF and these baseline cases in a dynamic user environment with load changes from moderate to low and then back to moderate, see Figure 14. The y-axis is the latency measured in ms, the x-axis represents the experiment's elapsed time. Fixed-low and Fixed-mod perform well under the loads that they are optimized for, but perform poorly when load changes. Fixed-other achieves more stable performance, but not best under any load levels. SERF outperforms all these baseline scheduling methods and consistently adapts to the load change to achieve lowest latency. This experiment validates the need for adaptivity of SERF in a dynamic workload environment, where, for example, a best configuration for high loads could be sub-optimal for low loads. In addition, the profiling cost of these fixed configurations is more than 2 orders of magnitude higher than SERF, e.g., for ImageNet-22K, it takes nearly 2 days to identify Fixed-low or Fixed-mod configuration by profiling, and it takes even longer for Fixed-other as profiling needs to be done for multiple loads. In comparison, SERF only takes a few minutes for identifying the best parallel configurations under various loads.

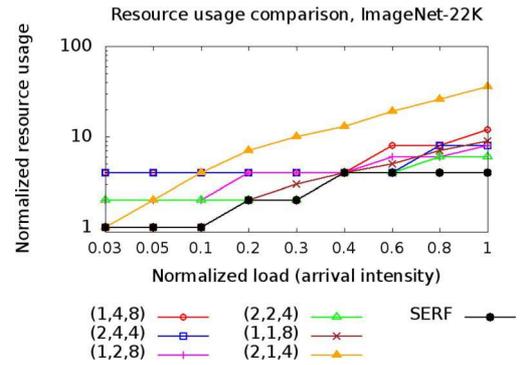


Fig. 15. Resource usage comparison between static configurations and SERF under different loads. The parameters of static configurations in the legend follows (max service parallelism, inter-node parallelism, intra-node parallelism).

G. Meeting SLO With Minimum Resources

To demonstrate the advantage of SERF over the static configurations selected based on heuristics, e.g., with minimum resource usage for certain loads, we compare their corresponding normalized resource usage required to meet the average latency SLO of 280 ms in Figure 15. Figure 15 suggests static configurations are best under certain loads, but become suboptimal for other loads while SERF always finds the best configurations that minimize the resource usage under different loads. In addition, the results also suggest that the resource usage gap between suboptimal configurations and best configuration is high, up to one order of magnitude, demonstrating the importance and effectiveness of using SERF.

H. Scalability and Applicability

When SERF works in large systems, the number of profiling experiments scales linearly with the total number of parallelism combinations. Because each profiling takes less than a few seconds, even for large systems running large and scalable applications with thousands of parallelism combinations, the profiling takes no more than a few hours. This profiling time can be further reduced to a few minutes if profiling experiments are conducted in parallel or in coarser granularity. In addition, the computation of the queueing model is efficient, i.e., constant with respect to the cluster size. Therefore, SERF is scalable to schedule large systems. We also stress that SERF is not limited to the CPU based system, it also works with systems using other architectures, e.g., GPUs, FPGAs, and ASICs. SERF can be used in different systems because it

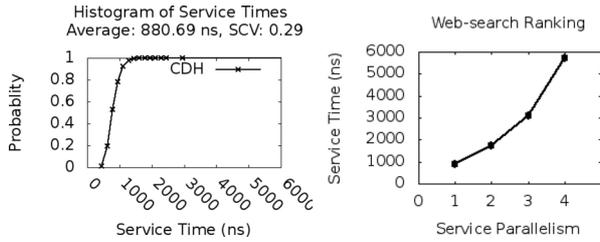


Fig. 16. CDH of service times (left) and service times under different number of concurrent running requests (right).

takes the basic system information as an input through simple profiling.

I. Generalization to Other Workloads

In this section, we demonstrate how SERF performs in general machine learning serving systems with similar workload characteristics (i.e., requests are homogeneous and there is interference among concurrent running requests). We use Web-search ranking [42] as an example for evaluation as it represents a typical supervised machine learning problem. In Web-search ranking, there is a query-document pair presented by a high-dimensional feature vector and the learning algorithm is to train a model to predict the relevance of a document to a query. In the serving stage, query-document pairs are requests sent to the serving system and the output is the relevance that is ranked based on the predictions from the trained model. We instrument the implementation in [42] to simulate a serving system. The original codes for computing the ranking is sequential, so we implement the service parallelism by using OpenMP. We did not implement the inter-node and intra-node parallelism as the request size is very small - the communication and synchronization overhead overcomes the benefits. We use Dataset 2 from Yahoo Ranking Challenge [43] as the workload. The testbed is a Dell PowerEdge R320 server with an Intel Xeon E5-2407 processor of 4 cores running at 2.2GHz and with 8 GB memory. The profiling results of service times indicate that the Web-search ranking serving system has similar workload characteristics as the DNN serving system, see Figure 16. The left plot in Figure 16 suggests that requests are homogeneous as the majority of requests are in the range of 700 ns to 1000 ns and their SCV is only 0.29. The right plot in Figure 16 shows very obvious interference among concurrent running requests as when the service parallelism increases, the average service time increases significantly.

We run extensive experiments with various load levels using different parallelism configurations and plot the average and distribution of the prediction errors in Figure 17. The results show SERF is quite accurate as the average error is only 1.07%, the 90th percentile is $< 3\%$, and the 95th percentile is $< 5\%$. We also show the prediction results of SERF for best parallelism configuration against the ground truth of measurement results in Figure 18. The results show that SERF always identifies the best configuration correctly.

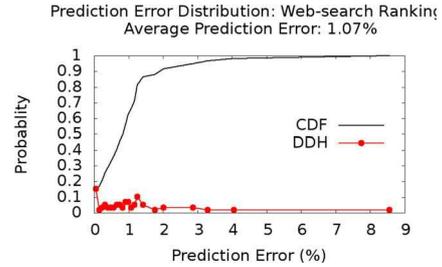


Fig. 17. CDF and DDH of prediction errors.

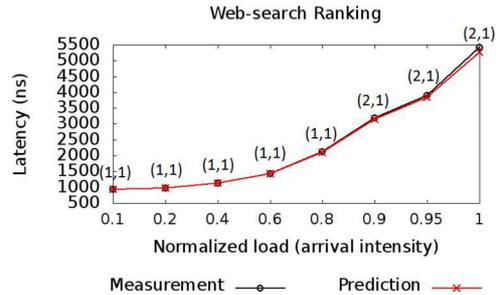


Fig. 18. Best configurations and according latency under different loads.

VI. RELATED WORK

DNN Serving: The state-of-the-art accuracy of DNNs on important artificial intelligence tasks, such as image recognition [1]–[3], speech recognition [9], [10], and text processing [11] has made the end-to-end latency of large-scale DNN serving systems an important research topic. Parallelism has been shown to be critical for good DNN performance at scale. Prior work [1], [3] has shown that parallel training on a cluster of commodity CPU machines achieves high throughput thus can train big DNN models (billions of connections) in a reasonable amount of time (days instead of months). Although these training platforms focus on improving system throughput instead of request latency, the parallelism mechanisms proposed there are directly translated to serving platforms as inter-node, intra-node and service parallelisms. Several recent work on DNN serving investigate hardware acceleration using GPUs [44], FPGAs [45], and ASICs [26]. They focus on mapping DNN computation to customized hardware, but parallelism has also been shown critical to offer low latency. For larger models that exceed on-chip RAM, [26] exploits inter-node parallelism across 64 ASIC chips with high-speed interconnect to achieve good performance. All these prior studies develop DNN serving platforms that support all or a subset of the parallelism mechanisms exploited in our paper. However, none of them investigates scheduling frameworks that make parallelism configuration choices based on DNN characteristics, hardware characteristics, and system load, which is the focus of SERF. Therefore, none of them quantifies the interference among concurrent running requests or various parallelism degree based on system load. Instead, they simply apply static partitioning of resources and use the maximum service parallelism available by the hardware [44]–[46]. SERF is complementary to the above work

and can be used as a scheduling framework for these serving platforms to identify best parallelism configurations and maximize their parallelism benefits.

While much of the recent work has focused on hardware acceleration using GPUs [20], [44], FPGAs [45], [46], and ASICs [26], [47], parallelism has been shown to be critical for good DNN serving performance at scale. For models that fit into on-chip RAM, service parallelism is exploited to improve serving throughput on FPGA [45], [46] and GPU [44]. For larger models that exceed on-chip RAM, the above specialized hardware that are significantly faster than CPU may not provide better performance. Reference [26] exploits inter-node parallelism across 64 ASIC chips with commodity high-speed interconnect to achieve good performance for larger models that exceed on-chip RAM. Compared to SERF, these prior work either do not consider the impact of request load and interference, or assume static partitioning of resources [44], therefore their performance and system utilization could benefit from SERF techniques.

Interactive Serving: There is a host of research in parallelizing request processing to reduce latency, and request scheduling in a multiprocessor environment to reduce latency. Here, we give a brief overview of systems that aim at quality of service and different parallelism options when resources are shared by many user requests. There has been a lot of work on measuring and mitigating interference among co-located workloads [48], [49]. The main theme is to predict performance interference among workloads and discover optimal workload co-locations to improve system utilization while meeting user performance goals. These studies treat each workload as a blackbox, and they do not consider solutions that involve modifying the workload (e.g., changing parallelism degree).

Adaptive allocation for server systems [50], [51] focuses on allocating resources dynamically to different components of the server, while executing each request sequentially, i.e., they consider service parallelism only. Adaptive parallelism for interactive server systems uses intra-node and service parallelism to reduce request latency. Recent work reduces request latency of interactive server systems using intra-node and service parallelism. Raman *et al.* [31] propose an API and runtime system for dynamic parallelism, where developers express parallelism options and goals, such as minimizing mean response time. Jeon *et al.* [30] propose a dynamic parallelization algorithm to decide the degree of request parallelism in order to reduce the average response time of Web search queries. Both approaches assume independent service time among requests, thus they do not consider interference among concurrent running requests, which is a key property of DNN workload supported by SERF. Another line of work [52] proposes to use parallelism to reduce tail latency based on the observation that requests exhibit large variability on service time, and therefore they use prediction [53] or dynamic parallelism [52] to execute long requests in parallel but short request sequentially. DNN requests, however, are homogeneous with similar service time, making these techniques ineffective. Moreover, none of these studies considers inter-node parallelism. Finding best parallel configurations has also been studied on other applications and systems, such as database, data analytics,

MapReduce [54]–[56]. However, none of these prior work leverages the distinctive properties of DNN workloads to exploit request homogeneity and interference awareness as SERF does.

Queueing Models: There is a vast area of research on queueing models, here we outline some Queueing models have been well studied, and here we outline some results that are related to the $M/D/c$ queue abstraction used in our work. While the solution of the $M/M/c$ system is exact [39], there are no exact solutions for $M/D/c$ systems. We note the existence of the Allen-Cunnen approximation formula for $GI/G/c$ [40] and Kimura’s approximation [57], both of which can also apply to $M/D/c$ since M is a special case of GI . Franx [58] provides an explicit expression for the waiting time distributions using probabilistic analysis while Cosmetatos [38] provides an approximate formula for $M/D/c$ that we use here. We enrich this approximation to account for the case where the service process is deterministic but also load dependent. Alternatively, an $M/D/c$ system can be approximated using an n -stage Erlang for the service process, essentially by approximating the system using a $M/Ph/1$ queue. While the $M/Ph/1$ queue can be solved using the matrix-geometric method [59], [60], the $M/Ph/c$ suffers from the well known problem of state space explosion. We direct the interested reader to [32] for an overview of various results on the $M/D/c$ queue that have been developed since the early 1930s. However, none of the above approximation methods for $M/D/c$ systems can be easily adapted to estimate latency of $M/D_{interf}/c$ systems. Here we extend the approximation by Cosmetatos to achieve this goal.

In [61], the basic idea of SERF that combines lightweight profiling with an interference-aware queueing-based prediction model is outlined and some initial evaluation results are presented. The current extended version of this paper provides a more detailed description of the SERF framework and presents a comprehensive performance evaluation study.

VII. CONCLUSION

In this work, we proposed a scheduling framework SERF for DNN serving systems that offers an automated way to optimally schedule serving requests, guaranteeing request serving that is done as fast as possible. A $M/D_{load}/c$ queueing system is at the heart of SERF: it is easily parameterized with limited workload profiling and it can provide a nearly instantaneous evaluation of many scheduling alternatives. SERF provides superior flexibility thanks to the following features: it provides adaptivity to load variations by quickly adapting scheduling parameters after changes in load are detected, it provides a mechanism to evaluate how changes in different parallelism can affect performance (and the related performance interference), and it quickly adapts to changes in user workload.

As service migrates to new clusters with different hardware, SERF offers again significant advantages to exhaustive profiling. Furthermore, thanks to the load dependent approach, the model can be easily adjusted for systems with heterogeneous

nodes: different $M/D_{load}/c$ models can be used for different hardware. Finally, SERF can support different scheduling objectives.

As future work, we plan to support more advanced control knobs, such as batching, in our framework. We also plan to make the framework accommodate orphan resource and straggler problem.

REFERENCES

- [1] J. Dean *et al.*, “Large scale distributed deep networks,” in *Proc. NIPS*, 2012, pp. 1223–1231.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Proc. NIPS*, 2012, pp. 1097–1105.
- [3] T. Chilimbi, J. Apacible, K. Kalyanaraman, and Y. Suzue, “Project Adam: Building an efficient and scalable deep learning training system,” in *Proc. OSDI*, Broomfield, CO, USA, 2014, pp. 571–582.
- [4] J. Mao, W. Xu, Y. Yang, J. Wang, and A. L. Yuille, “Explain images with multimodal recurrent neural networks,” *CoRR*, vol. abs/1410.1090, 2014. [Online]. Available: <http://arxiv.org/abs/1410.1090>
- [5] H. Fang *et al.*, “From captions to visual concepts and back,” in *Proc. CVPR*, Boston, MA, USA, 2015, pp. 1473–1482.
- [6] A. Karpathy *et al.*, “Large-scale video classification with convolutional neural networks,” in *Proc. CVPR*, Columbus, OH, USA, 2014, pp. 1725–1732.
- [7] J. Y. Ng *et al.*, “Beyond short snippets: Deep networks for video classification,” in *Proc. CVPR*, Boston, MA, USA, 2015, pp. 4694–4702.
- [8] S. Venugopalan *et al.*, “Translating videos to natural language using deep recurrent neural networks,” in *Proc. NAACL HLT*, Denver, CO, USA, 2015, pp. 1494–1504.
- [9] G. E. Dahl, D. Yu, L. Deng, and A. Acero, “Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 20, no. 1, pp. 30–42, Jan. 2012.
- [10] A. Y. Hannun *et al.*, “Deep speech: Scaling up end-to-end speech recognition,” *CoRR*, vol. abs/1412.5567, 2014. [Online]. Available: <http://arxiv.org/abs/1412.5567>
- [11] D. C. Ciresan, U. Meier, L. M. Gambardella, and J. Schmidhuber, “Deep big simple neural nets excel on handwritten digit recognition,” *CoRR*, vol. abs/1003.0358, 2010.
- [12] D. Talbot. *How Microsoft Cortana Improves Upon Siri and Google Now*. Accessed: Nov. 20, 2015. [Online]. Available: <http://www.tomshardware.com/news/microsoft-cortana-unique-features,26506.html>
- [13] R. Mcmillan. *How Skype Used AI to Build Its Amazing New Language Translator*. Accessed: Nov. 20, 2015. [Online]. Available: <http://www.wired.com/2014/12/skype-used-ai-build-amazing-new-language-translator/>
- [14] C. Rosenberg. *Improving Photo Search: A Step Across the Semantic Gap*. Accessed: Nov. 20, 2015. [Online]. Available: <http://googleresearch.blogspot.com/2013/06/improving-photo-search-step-across.html>
- [15] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. *A Picture Is Worth a Thousand (Coherent) Words: Building a Natural Description of Images*. Accessed: Nov. 20, 2015. [Online]. Available: <http://googleresearch.blogspot.com/2014/11/a-picture-is-worth-thousand-coherent.html>
- [16] B. Ramsundar *et al.*, “Massively multitask networks for drug discovery,” *CoRR*, vol. abs/1502.0207, 2015. [Online]. Available: <http://arxiv.org/abs/1502.0207>
- [17] F. Nelson. *Nvidia Demos a Car Computer Trained With Deep Learning*. Accessed: Nov. 20, 2015. [Online]. Available: <http://www.technologyreview.com/news/533936/>
- [18] A. Coates, B. Huval, T. Wang, D. J. Wu, and A. Y. Ng, “Deep learning with COTS HPC systems,” in *Proc. ICML*, Atlanta, GA, USA, 2013, pp. 1337–1345.
- [19] J. Deng *et al.*, “ImageNet: A large-scale hierarchical image database,” in *Proc. CVPR*, Miami, FL, USA, 2009, pp. 248–255.
- [20] Y. Jia *et al.*, “Caffe: Convolutional architecture for fast feature embedding,” in *Proc. ACM Int. Conf. Multimedia (MM)*, Orlando, FL, USA, Nov. 2014, pp. 675–678. [Online]. Available: <http://doi.acm.org/10.1145/2647868.2654889>, doi: 10.1145/2647868.2654889.
- [21] F. Bastien *et al.*, “Theano: New features and speed improvements,” *CoRR*, vol. abs/1211.5590, 2012. [Online]. Available: <http://arxiv.org/abs/1211.5590>
- [22] R. Collobert, K. Kavukcuoglu, and C. Farabet, “Torch7: A MATLAB-like environment for machine learning,” in *Proc. BigLearn*, 2011.
- [23] A. Krizhevsky, “Learning multiple layers of features from tiny images,” M.S. thesis, Comput. Sci. Dept., Univ. Toronto, Toronto, ON, Canada, 2009.
- [24] L. Bottou, “Large-scale machine learning with stochastic gradient descent,” in *Proc. COMPSTAT*, 2010, pp. 177–186.
- [25] N. P. Jouppi *et al.*, “In-datacenter performance analysis of a tensor processing unit,” in *Proc. 44th Annu. Int. Symp. Comput. Architect. (ISCA)*, Toronto, ON, Canada, 2017, pp. 1–12.
- [26] Y. Chen *et al.*, “DaDianNao: A machine-learning supercomputer,” in *Proc. MICRO*, Cambridge, U.K., 2014, pp. 609–622.
- [27] L. Dagum and R. Menon, “OpenMP: An industry standard API for shared-memory programming,” *IEEE Comput. Sci. Eng.*, vol. 5, no. 1, pp. 46–55, Jan./Mar. 1998.
- [28] A. D. Robison, *Intel Threading Building Blocks (TBB)*, 2011, pp. 955–964.
- [29] C. Demichelis and P. Chimento, “IP packet delay variation metric for IP performance metrics (IPPM),” IETF, Fremont, CA, USA, RFC 3393, 2002.
- [30] M. Jeon, Y. He, S. Elnikety, A. L. Cox, and S. Rixner, “Adaptive parallelism for web search,” in *Proc. EuroSys*, Prague, Czech Republic, 2013, pp. 155–168.
- [31] A. Raman, H. Kim, T. Oh, J. W. Lee, and D. I. August, “Parallelism orchestration using DoPE: The degree of parallelism executive,” in *Proc. PLDI*, San Jose, CA, USA, 2011, pp. 26–37.
- [32] H. Tijms, “New and old results for the M/D/c queue,” *AEU Int. J. Electron. Commun.*, vol. 60, no. 2, pp. 125–130, 2006.
- [33] C. Mega, T. Waizenegger, D. Lebutsch, S. Schleipen, and J. Barney, “Dynamic cloud service topology adaption for minimizing resources while meeting performance goals,” *IBM J. Res. Develop.*, vol. 58, nos. 2–3, pp. 1–10, Mar./May 2014.
- [34] Z. Zhang, L. Cherkasova, and B. T. Loo, “Optimizing cost and performance trade-offs for MapReduce job processing in the cloud,” in *Proc. NOMS*, Kraków, Poland, 2014, pp. 1–8.
- [35] J. Cao, W. S. Cleveland, D. Lin, and D. X. Sun, “On the nonstationarity of Internet traffic,” in *Proc. SIGMETRICS*, Cambridge, MA, USA, 2001, pp. 102–112.
- [36] M. F. Arlitt and C. L. Williamson, “Internet Web servers: Workload characterization and performance implications,” *IEEE/ACM Trans. Netw.*, vol. 5, no. 5, pp. 631–645, Oct. 1997.
- [37] T. Karagiannis, M. Molle, M. Faloutsos, and A. Broido, “A nonstationary Poisson view of Internet traffic,” in *Proc. INFOCOM*, Hong Kong, 2004, pp. 1558–1569.
- [38] G. P. Cosmetatos, “Notes approximate explicit formulae for the average queueing time in the processes (M/D/r) and (D/M/r),” *INFOR Inf. Syst. Oper. Res.*, vol. 13, no. 3, pp. 328–331, 1975.
- [39] L. M. Leemis and S. K. Park, *Discrete-Event Simulation: A First Course*. Upper Saddle River, NJ, USA: Pearson, 2006.
- [40] L. A. Baxter, “Probability, statistics, and queueing theory with computer sciences applications,” *Technometrics*, vol. 34, no. 2, pp. 240–241, 1992.
- [41] R. W. Wolff, “Poisson arrivals see time averages,” *Oper. Res.*, vol. 30, no. 2, pp. 223–231, 1982.
- [42] A. Mohan, Z. Chen, and K. Q. Weinberger, “Web-search ranking with initialized gradient boosted regression trees,” in *Proc. Yahoo Learn. Rank Challenge (ICML)*, 2011, pp. 77–89. [Online]. Available: <http://www.jmlr.org/proceedings/papers/v14/mohan11a.html>
- [43] Yahoo. *Yahoo! Learning to Rank Challenge*. Accessed: Nov. 20, 2015. [Online]. Available: <https://webscope.sandbox.yahoo.com/catalog.php?datatype=c>
- [44] J. Hauswald *et al.*, “DjiNN and Tonic: DNN as a service and its implications for future warehouse scale computers,” in *Proc. ISCA*, Portland, OR, USA, 2015, pp. 27–40.
- [45] C. Zhang *et al.*, “Optimizing FPGA-based accelerator design for deep convolutional neural networks,” in *Proc. FPGA*, 2015, pp. 161–170.
- [46] K. Ovtcharov *et al.* (2015). *Accelerating Deep Convolutional Neural Networks Using Specialized Hardware*. [Online]. Available: <http://research.microsoft.com/apps/pubs/default.aspx?id=240715>
- [47] T. Chen *et al.*, “Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning,” in *Proc. ASPLOS*, Salt Lake City, UT, USA, 2014, pp. 269–284.
- [48] J. Mars, L. Tang, R. Hundt, K. Skadron, and M. L. Souffia, “Bubble-up: Increasing utilization in modern warehouse scale computers via sensible co-locations,” in *Proc. ISCA*, Porto Alegre, Brazil, 2011, pp. 248–259.

- [49] R. Kettimuthu, G. Vardoyan, G. Agrawal, P. Sadayappan, and I. T. Foster, "An elegant sufficiency: Load-aware differentiated scheduling of data transfers," in *Proc. SC*, Austin, TX, USA, 2015, pp. 1–12.
- [50] G. Upadhyaya, V. S. Pai, and S. P. Midkiff, "Expressing and exploiting concurrency in networked applications with aspen," in *Proc. PPOPP*, San Jose, CA, USA, 2007, pp. 13–23.
- [51] M. Welsh, D. Culler, and E. Brewer, "SEDA: An architecture for well-conditioned, scalable Internet services," in *Proc. SOSP*, Banff, AB, Canada, 2001, pp. 230–243.
- [52] M. E. Haque *et al.*, "Few-to-many: Incremental parallelism for reducing tail latency in interactive services," in *Proc. ASPLOS*, Istanbul, Turkey, 2015, pp. 161–175.
- [53] M. Jeon *et al.*, "Predictive parallelization: Taming tail latencies in Web search," in *Proc. SIGIR*, Gold Coast, QLD, Australia, 2014, pp. 253–262.
- [54] X. Liu and B. Wu, "Scaanalyzer: A tool to identify memory scalability bottlenecks in parallel programs," in *Proc. SC*, Austin, TX, USA, 2015, pp. 1–12.
- [55] R. Susukita *et al.*, "Performance prediction of large-scale parallel system and application using macro-level simulation," in *Proc. SC*, Austin, TX, USA, 2008, pp. 1–9.
- [56] M. Sarwat, S. Elnikety, Y. He, and M. F. Mokbel, "Horton+: A distributed system for processing declarative reachability queries over partitioned graphs," *Proc. VLDB Endowment*, vol. 6, no. 14, pp. 1918–1929, 2013.
- [57] T. Kimura, "Approximating the mean waiting time in the GI/G/s queue," *J. Oper. Res. Soc.*, vol. 42, no. 11, pp. 959–970, 1991.
- [58] G. J. Franx, "A simple solution for the M/D/c waiting time distribution," *Oper. Res. Lett.*, vol. 29, no. 5, pp. 221–229, 2001.
- [59] M. F. Neuts, *Matrix-Geometric Solutions in Stochastic Models—An Algorithmic Approach*. New York, NY, USA: Dover, 1994.
- [60] A. Riska and E. Smirni, "ETAQA solutions for infinite Markov processes with repetitive structure," *Int. J. Comput.*, vol. 19, no. 2, pp. 215–228, 2007.
- [61] F. Yan, Y. He, O. Ruwase, and E. Smirni, "SERF: Efficient scheduling for fast deep neural network serving via judicious parallelism," in *Proc. Int. Conf. High Perform. Comput. Netw. Storage Anal. (SC)*, Salt Lake City, UT, USA, Nov. 2016, pp. 300–311.



Feng Yan received the B.S. degree in computer science from Northeastern University in 2008 and the M.S. and Ph.D. degree in computer science from the College of William and Mary in 2011 and 2016, respectively. He is an Assistant Professor of computer science and engineering with the University of Nevada Reno, Reno, NV, USA. His current research interests include big data, deep learning/machine learning, cloud computing, workload characterization, resource management, performance models and tools, scheduling, and storage systems. He is a member of ACM.



Yuxiong He received the B.Eng. degree in computer engineering from Nanyang Technological University in 2003 and the Ph.D. degree in computer science from Singapore-MIT Alliance, which is a joint graduate program of Massachusetts Institute of Technology, Nanyang Technological University, and National University of Singapore, in 2008. She is a Researcher with Microsoft Research, Redmond, WA, USA. Her current research interests include deep learning/machine learning, resource management, algorithms, modeling, and performance evaluation of parallel and distributed systems. She is a member of ACM.



Olatunji Ruwase received the B.S. degree in computer science from the University of Ibadan, the M.S. degree in computer science from Stanford University, and the Ph.D. degree in computer science from Carnegie Mellon University. He is a Senior RSDE with Microsoft Research, Redmond, WA, USA. His research interests include some combination of compiler, operating systems, and computer architecture techniques for understanding performance and reliability issues and crafting effective and practical solutions. He is a member of ACM.



Evgenia Smirni received the Diploma degree in computer science and informatics from the University of Patras, Greece, in 1987 and the Ph.D. degree in computer science from Vanderbilt University in 1995. She is the Sidney P. Chockley Professor of computer science with the College of William and Mary, Williamsburg, VA, USA. Her research interests include queuing networks, stochastic modeling, Markov chains, resource allocation policies, Internet and multi-tiered systems, storage systems, data centers and cloud computing, workload characterization, and models for performance prediction of distributed systems and applications. She has served as the Program Co-Chair of QEST'05, ACM Sigmetrics/Performance'06, HotMetrics'10, ICPE'17, and DSN'17. She also served as the General Co-Chair of QEST'10 and NSMC'10. She is an ACM Distinguished Scientist.