

# No DNN Left Behind: Improving Inference in the Cloud with Multi-Tenancy

Amit Samanta, Suhas Shrinivasan, Antoine Kaufmann, Jonathan Mace  
Max Planck Institute for Software Systems

## Abstract

With the rise of machine learning, inference on deep neural networks (DNNs) has become a core building block on the critical path for many cloud applications. Applications today rely on isolated ad-hoc deployments that force users to compromise on consistent latency, elasticity, or cost-efficiency, depending on workload characteristics. We propose to elevate DNN inference to be a first class cloud primitive provided by a shared multi-tenant system, akin to cloud storage, and cloud databases. A shared system enables cost-efficient operation with consistent performance across the full spectrum of workloads. We argue that DNN inference is an ideal candidate for a multi-tenant system because of its narrow and well-defined interface and predictable resource requirements.

## 1 Motivation

Deep neural networks (DNNs) excel at a wide range of machine learning tasks including computer vision, natural language processing, speech detection, and more. The success of DNNs has correspondingly led to the rapid growth of systems and platforms for deep learning (DL). Today, a rich ecosystem of platforms, libraries, and runtimes make it easy to develop, train, and deploy DNNs.

In a cloud and datacenter setting, machine learning workloads have thus grown in prominence. Broadly speaking, we can divide DL workloads into *training* workloads and *inference* workloads. Training is a compute-intensive batch task that constructs a DNN using large quantities of data; training bears similarity to other batch tasks like data analytics jobs and faces similar challenges. In contrast, *inference* is a low-latency, online task that generates predictions on-demand using a trained DNN; inference bears similarity to online applications like databases, web services, and microservices, and is often just one piece of a broader end-to-end application. DNNs are typically hosted separately from application logic and accessed via remote procedure call (RPC).

In this work we consider inference workloads exclusively. Although inference has not gone unnoticed in prior work, existing cloud infrastructure for inference workloads still has several limitations. These limitations are not fundamental, and we believe there are significant opportunities to better serve inference workloads. Throughout this paper we focus on deep neural networks (DNNs) as they are the main driver of today’s machine learning trends. However, many of our observations generalize beyond DNNs.

### 1.1 A Brief Primer on DNNs

For our purposes it suffices to think of a DNN as a stateless, deterministic, black-box function. For example resnet18, a DNN for image classification, takes as input a  $224 \times 224 \times 3$  RGB image, performs  $\sim 2$  billion flops, and outputs a 1000-dimension vector of class probabilities [22].

Internally, DNNs are straightforward, comprising a sequence of statically-defined *layers*. Each layer is a mathematical function that transforms the outputs from the previous layer to produce inputs for the next layer. For example, a *fully-connected layer* multiplies the output tensor of the previous layer with a tensor of hard-coded values (these were ‘learned’ during training). DNN depths vary from a few layers to a few hundred layers, and typically draw from a catalogue of several dozen layer types.

Evaluating a DNN thereby entails performing the operation of each layer in turn, transforming the original input into the output. DNNs can be evaluated on CPUs, GPUs, or special-purpose accelerators like TPUs [28]. GPUs and TPUs see significant speedups from parallelism; for example, for resnet18 with TVM [16] we measure 190.80ms median inference latency on a single-core CPU, compared to 0.97ms on a Tesla V100 GPU.

### 1.2 Deploying a DNN for inference

The conventional approach to deploying a DNN is to provision a container or virtual machine (VM), and within

the container or VM, host the model on a model server. Model servers are analogous to web servers, and notable examples include TensorFlow Serving [34], Apache MXNet Model Server [30], and Clippier [17].

Aside from dedicated VMs,<sup>1</sup> cloud customers can alternatively use a hosted system to serve their model, such as Google ML Engine [5] and Microsoft Azure ML [3]. These systems expose a higher-level interface for users – upload your model, and receive an endpoint to which you can send inference requests. Managed systems ease deployment complexity, as users do not have to manually provision specific resources or interact with underlying VMs. Internally, these systems also serve models using VMs. The main appeal of hosted solutions is to avoid capacity planning – these systems will automatically provision additional VMs if the workload increases, and alternatively tear down VMs if they are unused. Internal systems at companies such as Facebook, Google, and Uber take a similar approach [21, 23, 24, 34].

The common strand among these approaches is isolation at the VM level. Some model servers support deploying more than one model at a time, but this is statically configured on start-up, with pre-allocated resources [34]. All frameworks support models with custom user-code layers.

### 1.3 Expectations vs. Reality

All existing approaches to deploying DNNs use VMs for isolation. However, we argue that VMs are inefficient, undesirable, and fundamentally mismatched with the expectations and requirements of inference workloads. Inference workloads are online workloads and are often part of broader latency-sensitive end-to-end applications. The exact number of inference requests per second will vary unpredictably over time, often at fine time scales; meanwhile, the workload may have tail latency targets on the order of milliseconds, such as 7ms at Google [28], 10ms at Uber [23], and 25ms at Zendesk [15].

**Reality: idle resources and over-provisioning** With this in mind, the basic approach of statically provisioning VMs has clear drawbacks. First and foremost, users must either (i) over-provision to satisfy the estimated peak demand, thus leaving resources idle much of the time; or (ii) accept increased latency and even denial of service if demand increases. Moreover, some workloads may *never* have sufficient demand to warrant an entire VM, and may even undergo long periods of idleness. Statically-provisioned VMs are wasteful, as users must nonetheless provision and pay for the excess capacity. The problem is further compounded when using expensive hardware accelerators like GPUs and TPUs, which only make sense for workloads that can sustain thousands of requests per second [28]. Lastly, mapping workload requirements to

concrete resources is non-trivial. In a recent Amazon blog post, authors describe how “*developers are often stumped when the time comes to pick an instance type and size. Indeed, for larger models, the inference latency of CPUs may not meet the needs of online applications, while the cost of a full-fledged GPU may not be justified.*” [36]

**Reality: scaling in coarse-grained increments** Statically-provisioned VMs are further mismatched for workloads with tight latency requirements. To satisfy these latency requirements, it may be necessary to use hardware accelerators like GPUs and TPUs. Hardware accelerators can provide several orders of magnitude speedup for inference workloads; for example, for resnet18 we measure 190.80ms median inference latency on one-core Google Cloud VM, compared to 0.97ms on a Tesla v100. However, hardware acceleration is both more expensive and more *coarse-grained* than CPU-only VMs. Concretely, the throughput of resnet18 on one CPU core is 5.24 inferences per second (inf/s), at a cost of 3.48c/hr (\$1.84 per million inferences). Conversely the throughput of resnet18 on a GPU is 1031 inf/s.<sup>2</sup> GPUs are therefore more cost-effective at \$2.55/hr (\$0.69 per million inferences). However, there is clearly a tradeoff: though GPUs have the potential for better latency and throughput, CPUs are scalable in much finer increments. Each generation of GPU and CPU offers a different point on this scale.

**Reality: slow auto-scaling** Beyond static resources, cloud providers also offer managed solutions, such as Google ML Engine, and Microsoft Azure ML, outlined in §1.2. These solutions will scale provisioned resources in response to fluctuations in the workload. To do so, customers must specify latency targets or throughput thresholds, which trigger the system to automatically spin up new VMs or tear down VMs as workloads change. However, VMs are inherently slow to spin up, and as such, auto-scaling can only adapt to workload fluctuations over long time periods. For example, Amazon SageMaker [4] makes scaling decisions in 5-minute increments by default. Google Cloud ML’s documentation reflects this limitation: “*If your traffic regularly has steep spikes, and if reliably low latency is important to your application, you may want to consider manual scaling.*” [5].

**Reality: high cold-start latency** As a complement to auto-scaling, managed solutions also enable VMs to be torn down completely if workloads undergo long periods of idleness. This avoids paying for resources that will not be used, albeit over long time scales. However, when a workload picks up again, the resources must be re-provisioned. This task, sometimes referred to as *cold-start*, can take many seconds. For example, researchers evaluating the feasibility of DNNs in serverless applica-

<sup>1</sup> In the rest of the paper we use VMs to refer to both containers and VMs for brevity.

<sup>2</sup> With batching and batch-interleaving we can increase throughput up to 4083 inf/s at the cost of elevated 62.1ms median latency.

tions measured cold-start times of up to 12 seconds for 100MB models [26]. This increased latency is clearly at odds with the needs of online workloads, and in particular, any workload that may be infrequent or sporadic.

**Reality: mismatched pricing abstractions** Although managed solutions offer a high-level abstraction for inference workloads – “give us your model, then send us your inference requests” – this abstraction is not reflected in the pricing models of these systems. Both the underlying implementations of these managed solutions, and subsequent pricing models, are based on total VM time, *including idle time*, instead of per inference or total compute time, as is common in, *e.g.* serverless computing environments [9]. Consequently, the costs of idle resources due to over-provisioning are reflected back on the user. Given the long time scales over which these systems make decisions – *e.g.* waiting 10 minutes before tearing down a Cloud ML container [6] – this adds significant financial cost to users. This represents a fundamental tradeoff that is inherent to using VMs as the unit of provisioned resource: paying for unused, idle resources vs. incurring high cold-start latency.

## 1.4 Summary of Requirements

Based on these challenges, we summarize the following desirable properties for inference workloads:

**Latency** Inference workloads need stable average and tail latency, regardless of workload volume. Infrequent or sporadic workloads should not suffer from high cold-start latency.

**Elasticity** Online workloads are inherently unpredictable, and may rapidly transition from low-volume to high-volume. Workload fluctuations should be handled transparently.

**Cost-Effective** Users should not have to significantly over-provision resources. Moreover, costs should be consistent regardless of the workload volume, *i.e.* users should be charged based on the work done.

## 2 Multi-Tenant Systems

VMs and containers are a poor fit for inference workloads. Our proposed alternative is to share the underlying resources across tenants, by executing inference workloads within a *shared, multi-tenant system*. In shared systems, the system operator provisions resources for the system as a whole, and runs long-lived system processes that receive and execute requests from different tenants concurrently. By sharing resources, fluctuations in demand can be amortized across tenants, and we avoid over-provisioning and wasting resources. By sharing processes, workload spikes can be absorbed by re-distributing load, and workloads with long periods of

inactivity do not incur cold-start latency. And lastly, since cloud providers maintain control over the system and its resources, they can more closely align the pricing and system abstractions (*e.g.*, by charging on a per-request basis). Today, multi-tenant systems already provide a variety of core datacenter services, such as storage [13, 19], databases [1, 14], queueing [2], and co-ordination [12, 25].

Multi-tenant systems are only justifiable for *core data-center functionality*, where there is a common need for the functionality across many tenants and workloads. We believe that DNN inference is sufficiently important and prevalent to justify a specialized multi-tenant system. A shared system can significantly improve elasticity and cost-effectiveness. However, in designing a multi-tenant system for inference workloads, we must juggle our original goals with three additional challenges:

**Abstractions** Tenants no longer interact with compute resources directly. Instead, shared systems expose high-level abstractions and interfaces. These abstractions *must* generalize across many workloads and tenants.

**Security** Shared systems execute requests of different tenants within the same, shared processes. Thus, users are no longer separated by rigid OS or hypervisor boundaries. Nonetheless, we must still ensure security between different tenant workloads.

**Performance Isolation** The system must prevent performance interference between different tenants. Performance isolation is arguably the most difficult challenge for multi-tenant systems today, and remains an active area of research for multi-tenant systems in general [11, 31–33, 37].

In the remainder of this paper, we outline how a multi-tenant system for DNN inference can address these challenges, while also satisfying our original motivating requirements.

## 3 Design of a Shared Inference System

Our proposed system architecture bears similarity to many other shared systems [1, 14, 19]. High-level meta-operations are handled by a logically centralized controller. Meanwhile, DNN hosting and inference are handled by worker processes, spread across many machines.

Users first upload a DNN to the system and then start sending inference requests. Internally, the system distributes the DNN to one or more workers. Then, inference requests are routed to the workers that host the DNN. Workers host DNNs for many tenants simultaneously, and receive incoming requests. If a DNN is in high-demand, it is replicated across many workers. Replication achieves both fault tolerance, and elasticity. In our design, we focus particular attention on how workers multiplex inference requests across potentially many models.

### 3.1 Inference Runtime

Since the focus of our system is inference, a pre-requisite is for users to have trained their DNN using an existing deep learning framework such as TensorFlow [10]. However, today’s DL frameworks are designed to handle the entire model lifecycle, from training to deployment, and crucially, the framework also provides the runtime for model execution. In a multi-tenant system it is infeasible to allow users to upload an arbitrary runtime for executing their model, as this would require heavyweight isolation techniques to guarantee security, and it would limit opportunities for optimization. Instead, our system contains an internal model execution runtime and does not require user-code for performing inference.

However, using a custom runtime means users must submit models in a format understood by the runtime. Recent efforts in the machine learning community to address framework inter-operability have led to the development of the Open Neural Network eXchange (ONNX) model format [7]. The ONNX model format is simply a high-level description of the structure and parameters of a trained DNN, without framework-specific code or runtime optimizations.

The downside of this abstraction is that it restricts users to a pre-defined set of supported layer types. While existing DL frameworks allow users to implement custom layer types, user extensions are not feasible in a multi-tenant system. We do not believe that this is a significant limitation for our system targeting common-case workloads. For example, 95% of production DNNs at Google use standard layers; moreover, Google’s production hardware accelerator, the TPU, only supports standard layers [28]. Thus, supporting prototype layers is unnecessary until they reach mainstream adoption, at which point they can be included in our set of supported layer types. In general this restriction is reasonable and multi-tenant systems often disallow custom user code; for example, multi-tenant databases typically do not support SQL’s user-defined functions [1, 13].

### 3.2 Performance Isolation

In shared systems, performance isolation is important for ensuring aggressive tenants or unpredictable workloads do not cause starvation, reduced throughput, or high latency for others. However, comprehensive performance isolation is a challenge, even for existing systems today. Difficulties arise because isolation must be implemented at the application level, where we lack the ability to preempt requests while they are executing. A common approach is to predict resource requirements, measure actual consumption, and use coarse-grained feedback loops to provide guarantees like fairness over time [32, 35]. Often this is implemented as a fair queue scheduler at the request admission point [33].

However, unlike workloads from other domains, DNN inference has highly predictable resource consumption patterns. In other multi-tenant systems, performance isolation is difficult primarily because resource requirements are unpredictable and vary widely from request to request [33], and once a request is admitted it runs to completion. DNN inference does not face this challenge, because inference is a fundamentally predictable computation. This stems from the structure of DNNs (§1.1) – they are a fixed sequence of mathematical operations. A priori, we can quantify the exact number of flops required by each layer of the DNN. Moreover, DNNs are predictable as they do not contain control flow elements.<sup>3</sup> DNNs that accept variable-sized or batched inputs also vary deterministically based on input size.

We can exploit this predictability to do a much better job of scheduling, whether at request admission, or at finer granularity within the system. Instead of modeling costs up front, we propose a more pragmatic approach based on measurement. In our experiments with TVM [16] we measure 99th percentile latencies not exceeding 15% of the mean for a range of off-the-shelf DNNs [8] and workload mixes.

Predictable computations enable systems to react to workload fluctuations much more quickly, and enable higher quality scheduling decisions. For example, instead of heuristic-based best effort scheduling, an admission scheduler can confidently optimize an objective across all pending requests, such as minimizing average latency. Overall, predictable DNN inference presents an opportunity both to improve upon existing resource management techniques, and to explore new approaches entirely.

### 3.3 Efficiency and Optimizations

A key characteristic of a multi-tenant system is to alternate service between different tenants. In the worst case, each request may require loading and executing a different model that is not currently loaded. This introduces additional resource costs, such as the need to copy the model from a remote machine or from cold storage. Similarly, if we use hardware accelerators, then models need to be copied from host memory to device memory. Overall, the total inference latency will depend on a combination of execution latency (CPU, GPU, or other accelerator) and transfer latency (PCIe, disk, and/or network). As mentioned in §3.2, execution latency is predictable; but so too is transfer latency, since the memory footprint of a DNN is fixed. For example, resnet18 allocates approximately 78MB device memory for DNN weights; we measure ~7ms increased latency when copying weights prior to each inference, consistent with 12GB/s PCIe bandwidth.

Not all inference requests incur memory transfer over-

<sup>3</sup>This does not preclude higher-level control flow, which is the subject of recent research [27, 38]

heads because of caching opportunities at each level. The typical memory footprint of a DNN is in the tens or hundreds of MBs; in contrast, servers often exceed 1TB of main memory, current-generation GPUs have up to 32GB device memory, and current-generation TPUs have 64GB device memory. Most requests can exploit cached models instead of reloading from scratch. Consequently, bottlenecks will vary based on resource requirements of each DNN and cache hit ratios. For example, for models like resnet18, a GPU cache hit ratio of 85% or greater would shift the bottleneck resource from PCIe bandwidth to GPU execution latency.

In all, this leads to a multi-resource scheduling problem that, while similar to work in other domains [20, 32], has some unique constraints: (i) inference requests have predictable resource requirements; (ii) inference requests consume resources one-at-a-time; (iii) resources are independent and asynchronous; (iv) resources have measurable concurrency and throughput; and (v) scheduling decisions can be interposed before each resource. These constraints present an opportunity for high quality, fine-grained scheduling.

Beyond fine-grained scheduling decisions, we also have opportunities for high-quality placement and load-management decisions. For any request, we can calculate with high confidence the latency of local execution including any memory transfers. We can consider alternative execution strategies, such as CPU execution vs. hardware accelerator, and local vs. remote. A worker with several pending requests can calculate a priori the expected completion time of each request, including queuing time, and pre-emptively cancel or re-route requests accordingly.

Of course, the optimizations we have described primarily affect models with infrequent or varying workload patterns, for which multi-tenancy ‘reclaims’ resources that would otherwise go unused. This does not, however, come at the expense of degraded performance for heavy workloads. If a model does have a heavy workload (*i.e.*, enough demand to saturate a worker entirely), then we migrate colocated models elsewhere, giving the heavy workload essentially exclusive use of the worker. Then, any hardware [28] or software [18, 29] optimizations are equally applicable.

## 4 Discussion

Multi-tenancy has complementary goals to much of the prior work around DNN inference. Assuming sufficient workload demand from individual models, multi-tenant systems can equally benefit from specially designed accelerators [28], fine-grained batching techniques [18], and potential future results in inter-model batching [29]. A multi-tenant system would be particularly well-placed for exploiting inter-model optimizations, as the system

controls model placement and co-location decisions.

In the research literature, the most similar system to what we propose is Clipper [17]. Clipper also proposes an abstraction to serve as a “narrow waist” for model deployment, albeit different to our proposed abstraction, and isolates models using containers. Clipper focuses on challenges and optimizations that lie *above* their proposed interface – model management, latency-throughput trade-offs, and higher-level concerns like prediction accuracy. By contrast, we consider different challenges and optimizations that lie *below* our proposed interface. In industry, the most similar system to what we propose is TFS<sup>2</sup> [34], Google’s internal model hosting system, which distributes models to shared worker processes and provides automatic scaling; however, there is insufficient public information for a detailed comparison.

In this paper we did not discuss pre- and post-processing of DNN inputs and outputs, an important step for every DNN pipeline. We believe that this step is better handled by a separate (but possibly co-located and co-designed) system, that composes much like *e.g.* distributed file systems and databases. Processing steps have different performance characteristics compared to DNN inference, and often rely on user code; of course, this does not preclude entirely the possibility of safe high-level abstractions. The biggest difference between processing and inference is inter-model commonality. Pre- and post-processing steps are often similar between DNNs, and can be batched, even across different model pipelines [29]. However, DNN inference has fewer opportunities for batching across models, as model weights are unique.

Lastly, DNN inference does not cover all machine learning workloads. But, by restricting our design to this specific but common workload class, it enables assumptions around performance, predictability, and clarity, that we would otherwise lack. Models beyond DNNs have fundamentally different performance characteristics, thus we omit them from consideration. Similarly, we exclude reinforcement learning, which does not have a distinct inference phase. Multi-tenant systems for these other scenarios may also make sense, and we look forward to seeing future research in this direction.

## 5 Conclusion

In this paper we proposed that DNN inference should be a first-class cloud primitive, provided by a shared multi-tenant system. Multi-tenancy enables cost-efficient operation with consistent performance across a wide range of workloads. DNN inference is ideally suited for multi-tenancy because of its predictable resource requirements, which help address the important question of performance isolation.

## References

- [1] Apache HBase. Retrieved January 2019 from <https://hbase.apache.org>. (§2, 3, and 3.1).
- [2] Apache Kafka: A distributed streaming platform. Retrieved January 2019 from <https://kafka.apache.org>. (§2).
- [3] Deploy models with the Azure Machine Learning service. Retrieved January 2019 from <https://docs.microsoft.com/en-us/azure/machine-learning/service/how-to-deploy-and-where>. (§1.2).
- [4] Deploying a Model on Amazon SageMaker Hosting Services. Retrieved January 2019 from <https://docs.aws.amazon.com/sagemaker/latest/dg/how-it-works-hosting.html>. (§1.3).
- [5] Google Cloud Machine Learning Engine. Retrieved January 2019 from <https://cloud.google.com/ml-engine/docs/tensorflow/prediction-overview>. (§1.2 and 1.3).
- [6] Google Cloud ML Engine Pricing: More about prediction costs. Retrieved January 2019 from [https://cloud.google.com/ml-engine/docs/pricing#more\\_about\\_prediction\\_costs](https://cloud.google.com/ml-engine/docs/pricing#more_about_prediction_costs). (§1.3).
- [7] Open Neural Network Exchange Format: The new open ecosystem for interchangeable AI models. Retrieved January 2019 from <https://onnx.ai/>. (§3.1).
- [8] The ONNX Model Zoo. Retrieved January 2019 from <https://github.com/onnx/models>. (§3.2).
- [9] Microsoft Azure Functions. Retrieved January 2019 from <https://azure.microsoft.com/en-us/services/functions/>, 2019. (§1.3).
- [10] ABADI, M., BARHAM, P., CHEN, J., CHEN, Z., DAVIS, A., DEAN, J., DEVIN, M., GHEMAWAT, S., IRVING, G., ISARD, M., MANJUNATH, K., LEVENBERG, J., MONGA, R., MOORE, S., MURRAY, D. G., STEINER, B., TUCKER, P., VASUDEVAN, V., WARDEN, P., WICKE, M., YU, Y., AND ZHENG, X. TensorFlow: A System for Large-Scale Machine Learning. In *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2016), pp. 265–283. (§3.1).
- [11] ANGEL, S., BALLANI, H., KARAGIANNIS, T., O’SHEA, G., AND THERESKA, E. End-to-end Performance Isolation Through Virtual Datacenters. In *Proceedings of the 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2014), pp. 233–248. (§2).
- [12] BURROWS, M. The Chubby lock service for loosely-coupled distributed systems. In *Proceedings of the 7th USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2006), USENIX Association, pp. 335–350. (§2).
- [13] CALDER, B., WANG, J., OGUS, A., NILAKANTAN, N., SKJOLSVOLD, A., MCKELVIE, S., XU, Y., SRIVASTAV, S., WU, J., SIMITCI, H., ET AL. Windows Azure Storage: A Highly Available Cloud Storage Service with Strong Consistency. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP)* (2011). (§2 and 3.1).
- [14] CHANG, F., DEAN, J., GHEMAWAT, S., HSIEH, W. C., WALLACH, D. A., BURROWS, M., CHANDRA, T., FIKES, A., AND GRUBER, R. E. Bigtable: A distributed storage system for structured data. In *Proceedings of the 7th USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2006). (§2 and 3).
- [15] CHEE YAU, W. How Zendesk Serves TensorFlow Models in Production. <https://medium.com/zendesk-engineering/how-zendesk-serves-tensorflow-models-in-production-751ee22f0f4b>, February 2017. (§1.3).
- [16] CHEN, T., MOREAU, T., JIANG, Z., ZHENG, L., YAN, E., SHEN, H., COWAN, M., WANG, L., HU, Y., CEZE, L., ET AL. TVM: An Automated End-to-End Optimizing Compiler for Deep Learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2018), pp. 578–594. (§1.1 and 3.2).
- [17] CRANKSHAW, D., WANG, X., ZHOU, G., FRANKLIN, M. J., GONZALEZ, J. E., AND STOICA, I. Clipper: A Low-Latency Online Prediction Serving System. In *Proceedings of the 14th USENIX Symposium on Networked Systems Design and Implementation (NSDI)* (2017), pp. 613–627. (§1.2 and 4).
- [18] GAO, P., YU, L., WU, Y., AND LI, J. Low Latency RNN Inference with Cellular Batching. In *Proceedings of the 13th ACM European Conference on Computer Systems (EuroSys)* (2018), ACM, p. 31. (§3.3 and 4).
- [19] GHEMAWAT, S., GOBIOFF, H., AND LEUNG, S.-T. The Google File System. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP)* (2003), ACM, pp. 29–43. (§2 and 3).
- [20] GHODSI, A., SEKAR, V., ZAHARIA, M., AND STOICA, I. Multi-Resource Fair Queueing for Packet Processing. In *Proceedings of the 2012 Conference of the ACM Special Interest Group on Data Communication (SIGCOMM)* (2012). (§3.3).

- [21] HAZELWOOD, K., BIRD, S., BROOKS, D., CHINTALA, S., DIRIL, U., DZHULGAKOV, D., FAWZY, M., JIA, B., JIA, Y., KALRO, A., ET AL. Applied Machine Learning at Facebook: A Datacenter Infrastructure Perspective. In *Proceedings of the 24th IEEE International Symposium on High-Performance Computer Architecture (HPCA)* (2018), IEEE, pp. 620–629. (§1.2).
- [22] HE, K., ZHANG, X., REN, S., AND SUN, J. Deep Residual Learning for Image Recognition. In *Proceedings of the 29th IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 770–778. (§1.1).
- [23] HERMANN, J., AND DEL BALSIO, M. Meet Michelangelo: Uber’s Machine Learning Platform. <https://eng.uber.com/michelangelo/>, September 2017. (§1.2 and 1.3).
- [24] HERMANN, J., AND DEL BALSIO, M. Scaling Machine Learning at Uber with Michelangelo. <https://eng.uber.com/scaling-michelangelo/>, November 2018. (§1.2).
- [25] HUNT, P., KONAR, M., JUNQUEIRA, F. P., AND REED, B. ZooKeeper: Wait-free Coordination for Internet-scale Systems. In *Proceedings of the 2010 USENIX Annual Technical Conference (ATC)* (2010), USENIX Association. (§2).
- [26] ISHAKIAN, V., MUTHUSAMY, V., AND SLOMINSKI, A. Serving deep learning models in a serverless platform. In *Proceedings of the 2018 IEEE International Conference on Cloud Engineering (IC2E)* (2018), IEEE, pp. 257–262. (§1.3).
- [27] JEONG, E., JEONG, J. S., KIM, S., YU, G.-I., AND CHUN, B.-G. Improving the Expressiveness of Deep Learning Frameworks with Recursion. In *Proceedings of the 13th ACM European Conference on Computer Systems (EuroSys)* (2018), ACM. (§3).
- [28] JOUPPI, N. P., YOUNG, C., PATIL, N., PATTERSON, D., AGRAWAL, G., BAJWA, R., BATES, S., BHATIA, S., BODEN, N., BORCHERS, A., ET AL. In-Datacenter Performance Analysis of a Tensor Processing Unit. In *Proceedings of the 44th ACM/IEEE International Symposium on Computer Architecture (ISCA)* (2017), IEEE, pp. 1–12. (§1.1, 1.3, 3.1, 3.3, and 4).
- [29] LEE, Y., SCOLARI, A., CHUN, B.-G., SANTAMBROGIO, M. D., WEIMER, M., AND INTERLANDI, M. PRETZEL: Opening the Black Box of Machine Learning Prediction Serving Systems. In *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2018), pp. 611–626. (§3.3 and 4).
- [30] LUPESKO, H., WANG, Y., AND YU, R. Introducing Model Server for Apache MXNet. <https://aws.amazon.com/blogs/machine-learning/introducing-model-server-for-apache-mxnet/>, December 2017. (§1.2).
- [31] MACE, J., BODIK, P., FONSECA, R., AND MUSUVATHI, M. Towards General-Purpose Resource Management in Shared Cloud Services. In *Proceedings of the 10th USENIX Workshop on Hot Topics in System Dependability (HotDep)* (2014). (§2).
- [32] MACE, J., BODIK, P., FONSECA, R., AND MUSUVATHI, M. Retro: Targeted Resource Management in Multi-Tenant Distributed Systems. In *Proceedings of the 12th USENIX Symposium on Networked Systems Design and Implementation (NSDI ’15)* (2015). (§2, 3.2, and 3.3).
- [33] MACE, J., BODIK, P., MUSUVATHI, M., FONSECA, R., AND VARADARAJAN, K. 2DFQ: Two-Dimensional Fair Queuing for Multi-Tenant Cloud Services. In *Proceedings of the 2016 Conference of the ACM Special Interest Group on Data Communication (SIGCOMM)* (2016). (§2 and 3.2).
- [34] OLSTON, C., FIEDEL, N., GOROVY, K., HARMSSEN, J., LAO, L., LI, F., RAJASHEKHAR, V., RAMESH, S., AND SOYKE, J. TensorFlow-Serving: Flexible, High-Performance ML Serving. In *Proceedings of the Workshop on ML Systems at NIPS’17* (2017). (§1.2 and 4).
- [35] SHUE, D., FREEDMAN, M. J., AND SHAIKH, A. Performance Isolation and Fairness for Multi-Tenant Cloud Storage. In *Proceedings of the 10th USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2012). (§3.2).
- [36] SIMON, J. Amazon Elastic Inference – GPU-Powered Deep Learning Inference Acceleration. <https://aws.amazon.com/blogs/aws/amazon-elastic-inference-gpu-powered-deep-learning-inference-acceleration/>, November 2018. (§1.3).
- [37] YANG, S., LIU, J., ARPACI-DUSSEAU, A. C., AND ARPACI-DUSSEAU, R. H. Principled Schedulability Analysis for Distributed Storage Systems using Thread Architecture Models. In *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2018), USENIX Association, pp. 161–176. (§2).

- [38] YU, Y., ABADI, M., BARHAM, P., BREVDO, E., BURROWS, M., DAVIS, A., DEAN, J., GHEMAWAT, S., HARLEY, T., HAWKINS, P., ET AL. Dynamic Control Flow in Large-Scale Machine Learning. In *Proceedings of the 13th ACM European Conference on Computer Systems (EuroSys)* (2018), ACM. (§3).