

# Missing the Forest for the Trees: End-to-End AI Application Performance in Edge Data Centers

Daniel Richins<sup>1</sup>, Dharmisha Doshi<sup>2</sup>, Matthew Blackmore<sup>2</sup>, Aswathy Thulaseedharan Nair<sup>2</sup>, Neha Pathapati<sup>2</sup>, Ankit Patel<sup>2</sup>, Brainard Daguman<sup>2</sup>, Daniel Dobrijalowski<sup>2</sup>, Ramesh Illikkal<sup>2</sup>, Kevin Long<sup>2</sup>, David Zimmerman<sup>2</sup>, and Vijay Janapa Reddi<sup>1,3</sup>

<sup>1</sup>The University of Texas at Austin

<sup>2</sup>Intel

<sup>3</sup>Harvard University

## ABSTRACT

*Artificial intelligence and machine learning are experiencing widespread adoption in the industry, academia, and even public consciousness. This has been driven by the rapid advances in the applications and accuracy of AI through increasingly complex algorithms and models; this, in turn, has spurred research into developing specialized hardware AI accelerators. The rapid pace of the advances makes it easy to miss the forest for the trees: they are often developed and evaluated in a vacuum without considering the full application environment in which they must eventually operate.*

*In this paper, we deploy and characterize Face Recognition, an AI-centric edge video analytics application built using open source and widely adopted infrastructure and ML tools. We evaluate its holistic, end-to-end behavior in a production-size edge data center and reveal the “AI tax” for all the processing that is involved. Even though the application is built around state-of-the-art AI and ML algorithms, it relies heavily on pre- and post-processing code which must be executed on a general-purpose CPU. As AI-centric applications start to reap the acceleration promised by so many accelerators, we find they impose stresses on the underlying software infrastructure and the data center’s capabilities: storage and network bandwidth become major bottlenecks with increasing AI acceleration. By not having to serve a wide variety of applications, we show that a purpose-built edge data center can be designed to accommodate the stresses of accelerated AI at 15% lower TCO than one derived from homogeneous servers and infrastructure. We also discuss how our conclusions generalize beyond Face Recognition as many AI-centric applications at the edge rely upon the same underlying software and hardware infrastructure.*

## 1 Introduction

Artificial intelligence (AI), especially the field of machine learning (ML), is transforming the marketplace. Sparked by advances in computer system design, enterprises are leveraging AI in every possible manner to provide unprecedented new services to their end users, ranging from recommenda-

tion-based online shopping to personalized social network services, virtual personal assistants, and better health care.

To enable ML, there has been a flurry of work at two extremes. At one extreme is the effort that focuses on hardware acceleration of ML kernels [1, 2, 3, 4, 5]. At the other extreme is the effort that focuses on engineering the system and its supporting infrastructure, such as the associated networking and storage. The former is essential for enabling microprocessor advancements, while the latter is essential for allowing cloud-scale deployment.

But in recent years, we are seeing a shift in the needs of the industry. While much research has been dedicated to maximizing and accelerating machine learning performance, recent industry perspectives have urged for *a more holistic understanding of machine learning development and performance*. Facebook, for example, has discussed some of the challenges it has faced running AI at scale and encouraged research on mitigating those challenges [6]. Instead of focusing solely on the AI kernel computation time, there is a need to look at the bigger picture. Enabling AI applications involves several stages: ingesting the data, pre-processing the data, offloading the data to an AI accelerator, waiting for data, post-processing the result, etc., all of which affect the requests’ end-to-end latency and total system throughput.

At the same time, the industry is witnessing AI services migrate from warehouse-scale systems to smaller *purpose-built data centers* located at the edge, closer to end-users [7]. These edge data centers complement existing cloud- or large-scale services by being physically closer to the data source, which enables faster responses to latency-sensitive or bandwidth-hungry application services [8]. There are also data sovereignty and regulatory compliance rules to safeguard data privacy that are addressed with edge data centers [9]. Moreover, many mid-size organizations find it more economical to invest in on-premise data centers that are purpose-built for executing a particular type of task [10]. So, despite the continued growth in public cloud solutions, spending for edge data centers is predicted to increase [11].

In this work, we study the intersection of user-facing AI

computing and smaller, edge data centers to reveal the often overlooked “AI tax”: the additional compute cycles, infrastructure, and latency required to support the AI at the application’s heart. In the context of a data center, execution of a fully developed, deployment-ready AI-centric application *relies on more than just AI algorithms*. End users’ requests demand pre-processing to ready them for the pipeline; intermediate data must be communicated between stages, often over a network using custom protocols; the communication framework often has built-in data reliability safeguards which impose overheads on data movement; and each stage faces its own overhead for moving data. All of these components put together add to the overhead of executing AI.

We focus on video analytics in edge data centers for our studies of AI tax due to its rising importance. The global market for video analytics is expected to hit US\$25 billion by 2026 due to rapid adoption of video technologies across industries such as retail, manufacturing, and smart cities [12]. Video analytics uses AI to provide cost-efficient business intelligence insights to its users. The domain is slated for deployment in edge data centers, as opposed to traditional cloud- or warehouse-scale systems due to latency constraints, network bandwidth, and privacy regulations.

We study a full deployment of *Face Recognition*, an end-to-end video analytics AI-centric application at the edge. Our setup is an industry deployment of the Google FaceNet [13] architecture in an edge data center. As an AI-centric application, *Face Recognition* is a good choice as it employs three distinct artificial intelligence algorithms, including two neural networks and a classification algorithm. Furthermore, it is representative of the reality of a considerable portion of AI and ML applications: many AI applications exist as streaming services, deployed in data centers, serving real-time needs of consumers. Coordinating the many activities required to transform raw data into useful, easily consumable conclusions requires the intricacies and nuances of any distributed application: networking equipment, storage devices, coordination, data durability, power distribution, cooling, communication protocols, data compression, etc. [14].

We find that in today’s edge data centers, already *the communication framework can constitute over 33% on the latency of the application*. *Face Recognition* is built on top of Apache Kafka [15], which is widely adopted both directly and as a fabric upon which advanced streaming frameworks are built [16, 17, 18, 19, 20, 21]. Kafka is also representative of alternative frameworks that utilize communication hot spots. The simplicity and impressive performance of Apache Kafka have established it as a common denominator for many industry-quality projects. Despite this, requests can spend substantial time passing through the framework.

Moreover, we show that *as accelerator technologies advance and integrate into production environments, the supporting portions of the pipeline will soon supplant AI as the primary determiner of performance*. We measure the implications of greater AI inference acceleration. Apache Kafka becomes increasingly stressed to move the vastly increased volume of data ingested by the application. Even at relatively low acceleration factors, the added stress will quickly overwhelm Kafka’s current capabilities. We demonstrate that at a very modest  $8\times$  acceleration factor, Kafka

overwhelms the capabilities of its underlying storage.

These findings present a unique opportunity on the compute research spectrum: rather than neglecting the execution context of AI (missing the forest for the trees) and without moving into the realm of cloud compute where resources must be generic and homogeneous enough to handle all kinds of workloads, we show a proof-of-concept for the economic value of edge data centers. We demonstrate how a data center that is custom-built for the needs of a streaming AI workload can accommodate the anticipated requirements of accelerated AI without over-provisioning, thereby realizing an overall decrease in the total cost of ownership (TCO) in excess of 15% over a homogeneous edge data center.

Although our deep-dive primarily focuses on edge video analytics with *Face Recognition* (as the poster child application), our analyses and conclusions are not specific to that one application; we discuss how the underlying infrastructure of an end-to-end AI application will present mainly the same bottlenecks regardless of the AI application.

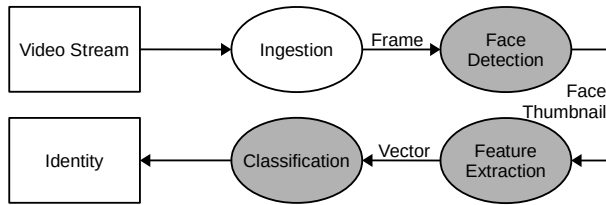
In summary, our main contributions and insights are

1. Where much focus is devoted to tuning and accelerating AI inference to enable faster compute, we instead evaluate the larger **system-level implications of end-to-end AI applications and expose the AI tax**;
2. We show that the **general-purpose CPU performance remains a significant determiner of overall request performance** because processing an end-user request requires more than just AI kernel computation;
3. The **communication layer of an AI application imposes a large overhead** on the latency of processing;
4. The **increased throughput from AI acceleration will overwhelm the communication substrate**; and
5. A **purpose-built data center can adapt to the upcoming challenges of accelerated AI at a lower TCO** than a generic, homogeneous data center.

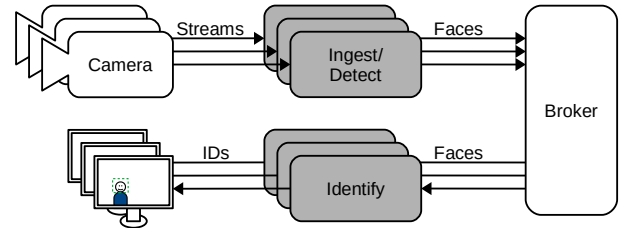
The remainder of this paper is structured as follows. In Section 2, we introduce our application, the running environment, and our experimental setup. In Section 3, we elucidate the AI tax, characterizing the performance and limitations of the end-to-end AI application. In Section 4, we conduct a forward-looking analysis of our application under accelerated AI compute and identify significant impediments to improving performance. In Section 5, we show that an edge data center can be purpose-built to address the upcoming challenges of AI while reducing TCO. In Section 6, we discuss opportunities for future work. We distinguish our work from prior art in Section 7 and conclude in Section 8.

## 2 AI in Edge Data Centers

The details of a typical edge application setup are generally proprietary. Hence, we contribute how we develop and deploy a user-facing video analytics pipeline, based on leading-edge AI algorithms, that runs in an edge data center. In Section 2.1 we introduce our video analytics AI workload, *Face Recognition*. We describe the implementation of *Face Recognition* in the context of a data center in Section 2.2. In Section 2.3, we present our edge data center hardware setup. We discuss the generality of our study and similar applications we have developed in Section 2.4.



(a) Algorithmic flow of *Face Recognition*. A video stream enters *ingestion* for separation into individual frames. *Face detection* finds any faces within a frame and produces a thumbnail for each. *Feature extraction* generates identifying features for each face. Finally, *classification* finds a nearest match to known faces to produce an identity.



(b) Data center deployment of *Face Recognition*. Algorithms run as standalone processes in lightweight Docker containers deployed on separate nodes throughout the data center. *Ingestion* and *face detection* both run in the *ingest/detect* container while *feature extraction* and *classification* are combined in *identify*. Communication between steps within a container happens internally while communication between containers is coordinated through Apache Kafka brokers.

Figure 1: From the conceptual operation of *Face Recognition* to a deployment-ready implementation. Shaded blocks show AI stages. The straight-forward application steps are wrapped in containers for deployment in a data center, requiring inputs, outputs, and communication coordination.

## 2.1 End-to-End Video Analytics Pipeline

Video analytics is the automatic analysis of video data. For our analysis, we developed and deployed a video analytics application for edge usage called *Face Recognition* (Figure 1). Though it uses machine learning, this application is strictly user-facing, i.e. it uses inference rather than training.

Our implementation of *Face Recognition* relies heavily on artificial intelligence and machine learning algorithms implemented in TensorFlow [22]. Given a number of input video streams, the application parses the videos into individual frames, locates faces within the frames, and identifies each face as a certain individual. The video streams could represent a surveillance system’s cameras [23], offline processing of recorded videos, a transactionless shopping environment [24], or many other applications where multiple streams are concurrently being fed into the system.

The *Face Recognition* application consists of four primary processing stages (Figure 1a) and is built from MT-CNNs (multi-task cascaded convolutional networks) [25] along with Google’s FaceNet [13, 26]. Like many real-world use cases, the application involves multiple inferences per query.

1. **Ingestion** is a pre-processing stage that ingests a video stream and parses it into individual frames. This stage is critical as FaceNet cannot operate directly on video.
2. **Face detection (AI)** relies on MT-CNNs to detect any faces within a frame without making any effort to identify them. It determines bounding boxes for and produces a 160x160 thumbnail of each face in a frame.
3. **Feature extraction (AI)** is built using the Inception-Resnet [27] architecture and produces a 128-byte vector of essential features that describe each face.
4. **Classification (AI)** compares the feature vector for a face against a set of known face vectors to find the best match by means of a support vector machine (SVM), yielding an identification.

*Face Detection* is designed as a pipelined streaming application—it ingests video streams at or near their native frame rate, injects them into the pipeline, and yields facial identities for frames after some delay. Even though the overall latency may exceed the time between adjacent frames in a video stream, because the application is pipelined, the throughput is at or near the native frame rate.

## 2.2 Edge Data Center Deployment

*Face Recognition* is more than just FaceNet—it is an application ready for deployment in an edge data center. Going from the algorithm to the full workload requires algorithm partitioning, containerization, work coordination, and communication management. We describe how the logical flow of *Face Recognition* is transformed into a functional streaming data center application (Figure 1b).

*Face Recognition* separates its algorithmic steps into discrete stages that coordinate with one another while running independently on separate nodes to produce a legitimate data center application. Separating the algorithm into multiple coordinated steps allows different stages of the application to adapt to the speed and requirements of other stages.

*Face Recognition* divides the algorithm into two stages (Figure 1b): *ingest/detect* and *identification*. Each is implemented as one or more processes running in a lightweight Docker [28] container. Each of the containers can be instantiated multiple times and deployed on one or more physical nodes in the data center. We also utilize a third container, the *broker*, that we discuss later.

*Ingest/detect* is a combination of ingestion and face detection. It runs an ingestion process, which accepts a video stream (we use a 1920x1080 video file for deterministic operation) and parses the stream into frames. It resizes the frames to 960x540 before passing them to the separate face detection process through an internal queue. Face detection produces a thumbnail for each face in a frame, if any (our video yields zero to five faces and averages 0.64 faces per frame, with face thumbnails averaging 37 kB each). If any faces are found, they must be transferred to the *identification* container, which is logically and physically separate from *ingest/detect*. Identification runs a single process, consisting of feature extraction and classification described in Section 2.1. It accepts faces from the *ingest/detect* container and internally produces a feature vector for each face which is compared against known faces to yield an identity.

In accordance with industry practice, we execute all inference directly on the CPU [6]. This yields the lowest latency which is critical in a user-facing application.

Communication between containers running on separate nodes within a data center requires intelligence and elegance.

We rely on Apache Kafka [15] to manage the communication between ingest/detect and identification, to allow for load balancing, and to dynamically adjust to node failures. Though there exists a variety of open source tools for building and managing streaming applications [29, 30, 31, 32, 33, 34], we note that these tools tend to rely on a separate framework for enabling communication between containers. It is common in practice to rely on Apache Kafka to serve this purpose, and each of these projects has proponents extolling the benefits of using Kafka [16, 17, 18, 19, 20, 21]. We therefore use Kafka directly to coordinate communication between our containers.

Apache Kafka implements the publish-subscribe pattern of communication [35]. This pattern operates by relying on an intermediate staging area for data, instead of data producers sending data directly to data consumers. Data *producers* publish data without any knowledge of the data *consumers* or even a guarantee that any consumers exist. They simply publish the data as a *topic*—a category of data where each individual item is treated identically. Similarly, consumers subscribe to a topic, oblivious to all details about the producers. As producers publish data to a given topic, the data become visible to the consumers subscribed to the same topic; consumers are then free to process the data.

In Kafka, the intermediate staging area where topic data is stored is implemented in *brokers*. A topic is implemented by creating partitions—open file handles—typically spread across multiple brokers. When a producer publishes data to a topic, it may send that data to any of the partitions, which the corresponding broker receives and writes to the open file. When a consumer requests data from a partition, the broker reads it from the same file. In contrast to producers, partitions may have a maximum of one consumer. Thus an application should divide a topic into at least as many partitions as there are consumers in order to maximize parallelism.

The topic partition also serves as the basic unit of replication. Kafka assumes and encourages data replication for reliability should a broker go down. Each partition has a “leader” and, in the presence of replication, some number of followers. After new data is written to a leader partition it is replicated to the followers. Producers and consumers interact with the broker that holds the leader partition, while the follower partitions are spread among the remaining brokers. In the event of a broker failure, one of the follower partitions will become the new leader partition. Unlike partitions, there are no “leader brokers” or “follower brokers.” Both leader and follower partitions are spread among all available brokers; thus, no one broker is more important or heavily utilized than any other.

For *Face Recognition*, the ingest/detect containers function as producers, sending face thumbnails extracted from each frame to brokers as the “faces” topic. The identification containers are the corresponding consumers, subscribing to the “faces” topic. The placement of brokers between ingest/detect and identification containers was chosen to provide load balancing. As we will show in Section 3, the two containers have different latencies; we thus instantiate more identification than ingest/detect containers. By placing brokers between them, Kafka ensures that the work is spread among the consumers evenly.

**Table 1: Server details. Each server in our data center is well-equipped, using leading-edge technology. Our nodes are powered by Intel Xeon Platinum 8176 or comparable CPUs.**

| Component           | Details                          |
|---------------------|----------------------------------|
| CPU                 | 2× Intel Xeon Platinum 8176 [37] |
| Cores               | 28                               |
| Base Frequency      | 2.10 GHz                         |
| Max Turbo Frequency | 3.80 GHz                         |
| SMT                 | 2-way                            |
| LLC                 | 38.5 MB                          |
| Memory              | 384 GB DDR4-2666 [38]            |
| Storage             | Intel SSD P4510                  |
| Read BW             | 2.85 GB/s                        |
| Write BW            | 1.1 GB/s                         |
| Read Latency        | 77 us                            |
| Write Latency       | 18 us                            |
| Network             | Full duplex 100 Gbps Ethernet    |

An alternative setup would separate the ingest/detect container into two separate containers, relying on the brokers to transfer data between them through an additional “frames” topic. While we do not show the data in this paper, we explored this option extensively, ultimately rejecting it in order to reduce network traffic (the importance of which will become apparent in Section 4).

## 2.3 Edge Data Center Configuration

We utilize a small edge data center built from high-performance servers (see Table 1). We use over 2200 processor cores spread across 40+ nodes to ensure that we have a realistic deployment whose characteristics can scale to larger setups. Each node is equipped with 56 physical cores spread across two sockets, 384 GB of RAM, high-speed local NVMe storage, and 100 Gbps Ethernet. The nodes are connected in a fat tree topology [36].

We rely on Kubernetes [39] and Docker [28] to deploy our application containers throughout the data center, reflecting common industry practice. Each of the three distinct pre-built container images (broker, ingest/detect, and identification) is deployed a set number of times and distributed throughout the data center. Because of the extremely low network utilization relative to capacity (Section 4.4), the placement of containers relative to one another in the data center is unimportant. We rely on Kubernetes to manage the deployment of containers. We use a minimum of three broker nodes in all cases to allow for three-way data replication, reflecting common practice in industry-quality deployments.

We omit the final stage represented in Figure 1b (displaying the identities on the video feed) from measurement studies because it offers little to no real computational insights.

## 2.4 Generalizability and Other Workloads

Beyond the value of an industrial case study,<sup>1</sup> we have also developed similar edge deployments of two additional applications—*Object Detection* and *Language Translation*—also built using the same infrastructure, including Kafka. While we do not study these in the same detail as *Face Recogni-*

<sup>1</sup>Case studies are an important but underappreciated tool in computer systems research [40]. They shed light on an application that is valuable in its own right and can pioneer an evaluation approach. The results from a case study are also often generalizable.

tion, we present some results in Section 3.4 showing that they exhibit the same kinds of behaviors as *Face Recognition*. *Object Detection* uses an R-CNN [41] to uniquely identify multiple objects in a frame. *Language Translation* takes advantage of Google’s neural machine translation solution (GNMT) [42] to translate text between English and German.

### 3 AI Tax

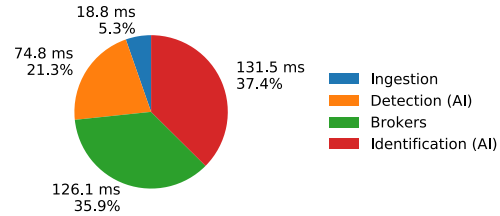
We start our exposition of the AI tax by evaluating the end-to-end performance of *Face Recognition*. We aim to understand what fraction of the cycles in an AI application go to AI processing versus the non-AI components. To this end, we examine the lifetime of a frame as it flows through the AI-centric *Face Recognition* application. In Section 3.1, we explain how we measure frame progress. Section 3.2 breaks down the end-to-end progress of a frame in each stage of the pipeline and shows that AI computation is not so central as one would expect in an AI application. In Section 3.3 we break down the application behavior in each container and reveal how much supporting compute is needed to enable AI processing. We show that it is vitally important to view AI application performance holistically, as it involves much more than just AI processing and the supporting code and infrastructure tax have a profound impact on latency. Finally, in Section 3.4, we present some basic data from other applications suggesting that they will face the same issues.

#### 3.1 Instrumenting the End-to-End Execution

To really understand an AI application deployed in even an edge data center, we must raise the level of abstraction from how applications are traditionally evaluated. While we do not claim to have the right level of abstraction for all end-to-end workloads, for *Face Recognition*, we believe that we have identified a good level of abstraction for tracking and measuring application progress without perturbing the application’s original behavior.

Application progress is a sequence of unit steps that are necessary for a frame to progress through the application. We term the units of application progress “events”; these are high-level steps in the application and roughly correspond to the stages introduced in Section 2.1: video ingestion, face detection, broker waiting time, and identification. Event-based logging lets us track end-to-end application progress. This higher level of abstraction is critical in enabling engineers to architect at a cluster level, where the complete application executes, instead of just at a node level.

We log all the events during execution of the application using Elasticsearch [43] and Logstash [44] running on a separate server. We measure the execution time of each step as well as the sizes of data that are transferred between stages. This is done using timestamps around the major regions of interest, e.g. the time to do the face identification *excluding* the supporting code (e.g. iteration management). In essence, events capture the major steps that a frame goes through from beginning to end. We use built-in language functions to measure the size of the data that are transferred through the brokers. Due to the infeasibility of instrumenting a complex program such as Kafka, we approximate the broker waiting time event by calculating the time delta between the end of face detection and the beginning of identification.



**Figure 2: Breakdown of end-to-end frame latency. With inference steps in detection and identification stages, less than 60% of the latency arises from AI stages. Over a third of a frame’s lifetime is spent in brokers.**

Our instrumentation method has negligible overhead and resource requirements, since we are only logging events. It has minimal impact on the application (see Figure 3).

#### 3.2 AI Applications Are More Than Just AI

One of the end user’s primary concerns is latency. In *Face Recognition*, this is the total time of a frame progressing serially from ingestion through identification; the latency of any stage that is not performing AI contributes to the AI tax.

We conduct experiments with 840 ingest/detect processes (producers) executing on 15 nodes (56 processes per node), 1680 identification processes (consumers) executing on 30 nodes (56 processes per node), and 3 brokers (each given its own node). We require the brokers to maintain 3× data replication, which is standard practice for disaster recovery.

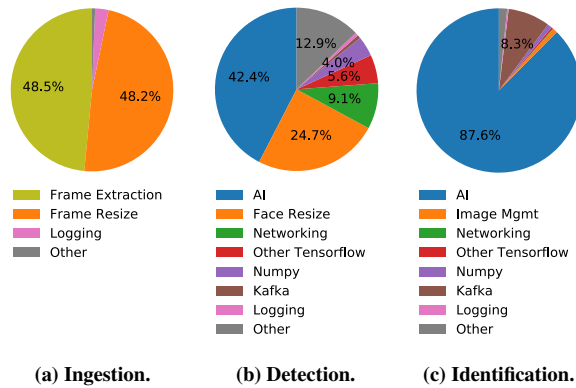
We measure an average face size of 37.3 kB and an end-to-end latency of 351 ms. While this latency may seem large, there are two points to remember. First, the throughput per stream is around 10 frames per second (FPS)—and a single stream could be divided among three ingest/detect instances for 30 FPS operation—regardless of the latency, since the application is pipelined; the output video still displays smoothly, just with a small delay. Second, there are multiple inferences per frame, performed sequentially, with the inference stages located on different nodes to improve performance [46]; the communication between the stages imposes additional latency.

Figure 2 summarizes the average latency for each stage. Ingestion operates quickly, taking only 18.8 ms, while the AI stages, face detection and identification, take 74.8 and 131.5 ms, respectively. Remarkably, over a third of the end-to-end latency is spent waiting between stages, at 126.1 ms. As in any real-time application, tail latency is an important factor to consider. We measure a 99<sup>th</sup> percentile tail latency of 2.21 s, with the standalone 99<sup>th</sup> percentile tail latencies of ingestion, detection, waiting time, and identification at 27 ms, 1.84 s, 116 ms, and 380 ms, respectively.

The end-to-end tail latency derives almost entirely from the waiting time in the brokers, which results from congestion in the application. When ingest/detect produces a surplus of faces, identification has a hard time keeping up, leaving the faces in the brokers for a longer time. When there are almost no faces detected, identification containers are almost idle and so are able to fetch the few faces quickly.

In summary, deep learning inference performance is more than just the performance of an individual node in the system. Even with a well balanced system, there is a substantial AI tax latency imposed in managing the transfer of data be-





**Figure 3: Process CPU time breakdowns.** Though ingestion uses no AI algorithms, it consists of straightforward image processing algorithms, which should be easily accelerated [45]. Face identification is overwhelmingly AI-centric. In contrast, face detection relies heavily on supporting code to enable its AI processing.

tween the nodes (i.e. the detection and identification stages). Without looking at the end-to-end latency, one would not realize that a large portion of time is spent waiting at brokers. This observation is *not* unique to our application; any application built on Apache Kafka (or a similarly brokered communication mechanism) will face this reality.

### 3.3 Overhead of Pre- and Post-Processing AI

Most AI research papers focus only on the core AI component, neglecting the other associated parts that are essential to end-to-end AI processing. However, there are pre- and post-processing steps that are unavoidable. Both play a critical role in the overall latency and both contribute to the AI tax. Pre-processing involves preparing the data for the AI kernel execution, while post-processing is loosely defined as any processing that is performed to convert the generated AI result(s) into something meaningful to the user or next stage.

To quantify the AI tax for pre- and post-processing we refine our view of a frame’s processing by looking at the time breakdown of each process using code profiling tools. Figure 3 shows, for each of the main processes (ingestion, detection, and identification), where time is spent.

Ingestion is exclusively a pre-processing stage. It shows a nearly even split between frame extraction and frame resizing (Figure 3a). Extraction refers to parsing the incoming video stream into individual frames. Resizing converts frames from 1920x1080 to 960x540 for the detection stage. The remaining time is split between the overhead of event logging and other supporting code, including transferring frames to the co-located detection process.

During face detection (Figure 3b), despite being an AI-centric stage, only 42% of the time is spent executing the AI algorithm in TensorFlow. Cropping and resizing faces (to 160x160) for identification takes 25% of the time; supporting TensorFlow and Numpy code (pre- and post-processing for each frame) take 6% and 4%, respectively; and “other” code takes a whopping 13% of the time. Code in the “other” category includes inter-process communication (from the ingest stage), additional matrix manipulation, loop management, bounding box calculation, image encoding, etc.

The AI-centric identification stage has a markedly differ-

ent breakdown. It spends 88% of its time directly executing AI algorithms; Kafka code, though, takes 8% of the time. The remaining components contribute little to the total time.

Beyond the pre-processing of the ingestion stage, end-to-end *Face Recognition* requires substantial pre- and post-processing within the AI-centric stages. In face detection, non-AI computation constitutes 57.6% of the compute cycles. In identification, that figure drops to 12.4%, which is still far from trivial. In a complex and diverse AI-centric application such as *Face Recognition*, AI computation constitutes 55.2% of end-to-end cycles, with the remainder going to supporting code: 17.8% to resizing, 9.0% to networking, 5.2% to tensor preparation, 3.6% to Kafka processing (outside of the brokers), and the rest to other supporting tasks.

In summary, despite the massive excitement surrounding AI algorithms, AI workloads are more than just tensors and neural networks: without the supporting code, AI is impotent. We emphasize that the supporting code, far from being a minor player in a complete application, constitutes over 40% of the compute cycles, not counting the compute time in the brokers. The pre- and post-processing code is executed on the general-purpose CPU, so it motivates the need to understand the role of the CPU as AI acceleration increases.

### 3.4 Additional Applications

Though we do not study them in as great detail as *Face Recognition*, our other applications show behaviors that indicate they will face the same issues as *Face Recognition*.

*Object Detection*, though also built from two stages and using Kafka as a communication substrate, only utilizes AI in one stage. This inference stage is computationally intensive; at 475.2 ms average compute time, its cycles are 97.2% ML. However, instead of transferring faces or objects through the brokers, it passes entire video frames. These weigh in at an average of 86.2 kB; image compression keeps this from ballooning too large. The significance of the larger data size will become apparent in Section 4.4.

*Language Translation* similarly has only one AI stage, which spends 79.4% of its cycles on ML. The large fraction of supporting code reemphasizes the importance of general CPU compute. Since the data passed through the brokers is purely textual sentences, the data transferred is small, typically no more than a couple hundred bytes per sentence.

## 4 Accelerating AI

There are numerous efforts underway to accelerate AI [1, 2, 4, 47, 48, 49, 50]. But given the significance of the AI tax in end-to-end AI performance and in pre- and post-processing, it is important to understand how the tax evolves as AI is accelerated and its impact on the overall end-to-end application performance; there are performance limitations that arise beyond a certain point of AI acceleration. To build a balanced system, it is important to understand these limits. We therefore study how varying degrees of AI speedup affect the end-to-end performance of our video analytics workload.

We do not attempt to promote any one or group of acceleration techniques, nor do we presume that a particular approach will be more successful than any other. In general, we could foresee accelerated AI coming about in various ways: CPU manufacturers may decide to integrate ML-

centric hardware directly into the CPU execution pipeline; or dedicated off-chip accelerators may be utilized, including highly-parallel pipelines (such as GPUs) and dedicated inference engines (such as Intel’s Neural Compute Stick [4], Habana’s Goya inference processor [51], or Google’s Coral Edge TPU [50]). Comparisons of the efficacy of each of these solutions is the subject of other work. In this work, we look at the impact of theoretical speedups, regardless of how the speedups are achieved.

In this section, we explore the impact of accelerating AI applications up to  $32\times$ , based on the performance of existing accelerators. Habana reports that its Goya processor achieves  $13.5\times$  speedup over a two-socket Intel Xeon Platinum 8180 [52].<sup>2</sup> We explore speedups approximately twice as great as this in order to account for future advances.

Section 4.1 analytically estimates accelerated AI performance to show its asymptotic limits. Section 4.2 introduces our technique for emulating accelerated workloads on current hardware. In Section 4.3 we evaluate the performance of accelerated processing and discover a quickly approaching bottleneck. Section 4.4 shows that the bottleneck results from overwhelming the system’s capacity to write to storage. We finish by showing in Section 4.5 how frames’ waiting time in brokers grows as a fraction of end-to-end latency.

#### 4.1 Analytical Speedups for AI Acceleration

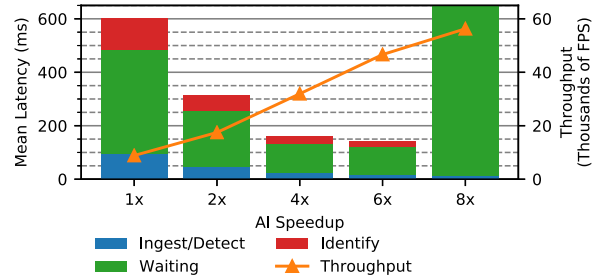
The AI tax means that a significant portion of an AI application’s compute cycles are spent on tasks other than AI and ML, and Amdahl’s law dictates that the overall speedup of a system is limited by the portion of execution that is not accelerated. Application of Amdahl’s law shows that each of the three primary processes—ingestion, detection, and identification—is limited in how much real-world speedup it can enjoy if AI is accelerated in isolation. Ingestion, which performs no AI compute, naturally derives no benefit from acceleration. Detection, which is 42% AI, rapidly approaches its asymptotic speedup of just  $1.74\times$ , achieving  $1.59\times$  overall speedup at  $8\times$  acceleration and  $1.66\times$  overall speedup at  $16\times$  acceleration. Identification, at 88% AI, has an asymptotic speedup limit of just  $8\times$ . At  $16\times$  AI acceleration it achieves  $5.6\times$  overall speedup, and even at  $32\times$  AI acceleration it shows just  $6.6\times$  overall speedup.

The exciting speedups promised by up-and-coming inference accelerators will be severely moderated by the reality of the supporting, non-AI code—the AI tax. With the fervor surrounding acceleration of AI and ML, these results from Amdahl’s law serve as an important reminder that AI applications are more than ML computation. Supporting and enabling code is a critical component of an end-to-end application and this should serve as a call to action to address the limitations imposed by that code.

#### 4.2 Emulating AI Acceleration on Hardware

It is instructive to see how the AI tax evolves as compute is universally accelerated (i.e. overcoming the asymptotic limits of Section 4.1) on a real system. To do so, we emulate the behavior of accelerated processing. Only the most basic loop controls and Kafka code are left in their original state.

<sup>2</sup>Other than its faster clock, the 8180 CPU is identical to the 8176 CPUs we use.



**Figure 4: Average frame latency and throughput under increasing AI acceleration. Beyond  $8\times$  speedup, the increased throughput leads to an unbalanced system. Queuing theory dictates that if elements enter the system faster than they leave, the latency increases to infinity.**

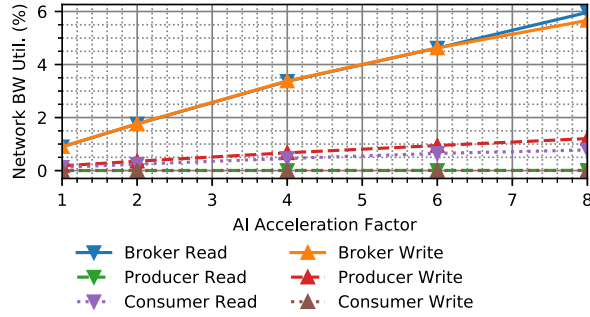
Our emulated acceleration technique relies on the observation that, from the perspective of application progress, the perspective of network traffic, and the perspective of the brokers, it is impossible to distinguish between (1) running the real application as has been described and characterized and (2) implementing artificial delays reflective of the actual compute times (Section 3) and sending meaningless data over the network of the same size as in the real application. In accelerated *Face Recognition*, rather than accelerating and executing the real algorithms, we replace the compute with calls to `sleep`, where the sleep duration is reflective of measured execution times (Section 3.1). Accordingly, rather than sending face thumbnails to brokers, we send meaningless data whose size matches the measured sizes. We can accelerate processing (both ingest/detect and identification) by an arbitrary factor by dividing the sleep times by the speedup factor. In this way, we maintain the behavior of the brokers, network, storage, and supporting code while exploring how acceleration changes the AI tax.

We emphasize that this AI acceleration emulation provides realistic performance estimation under acceleration because (1) the most basic general purpose processing (the support code to iterate through available frames, code to coordinate communication with Kafka brokers, the brokers themselves, etc.) remains in place and is executed as usual without the benefits of acceleration; (2) from the perspective of the data center, compute time spent executing real algorithms and waiting in `sleep` are identical; and (3) the brokers are completely ignorant of and unconcerned with the execution details of both producers and consumers. Thus, our setup to accelerate AI through emulation provides a realistic and believable look at the impact of faster AI on the data center and on the workload as a whole.

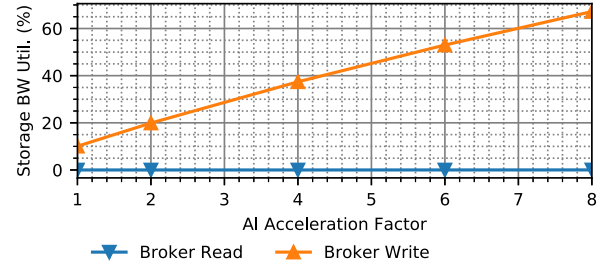
#### 4.3 Accelerated AI Impact on Total Speedup

We explore how the end-to-end frame latency will evolve as AI benefits from increasingly powerful acceleration. For this analysis, we assume that the AI algorithms will experience no latency overhead from acceleration—that is, we assume that the latency to communicate with dedicated off-chip accelerators is factored into the emulated speeds, or, equivalently, that future CPU architectures will directly integrate accelerator hardware into their execution pipelines.

For these experiments, we maintain the same application



(a) Network bandwidth utilization. Network activity from and to producers and consumers is concentrated at the brokers; nevertheless, broker network utilization falls far short of our 100 Gbps capacity.



(b) Storage bandwidth utilization. Storage activity for producers and consumers is not shown because it is not preserved in accelerator emulation. Network activity is translated to storage activity in the brokers.

Figure 5: Network and storage bandwidth utilization under acceleration. Whereas network utilization never exceeds 6% of network capacity, storage bandwidth utilization exceeds 67% of capacity at 8x acceleration. Storage bandwidth becomes a bottleneck much sooner than network bandwidth.

organization as depicted in Figure 1b. As we accelerate producers and consumers equally, the imbalance between the two will persist, so it is important to utilize Kafka’s load balancing features at the interface between the two.

For the sake of simplicity and repeatability and without loss of generality, we configure these emulation experiments so that each frame produces exactly one face. This has two impacts on performance: (1) on average, this is more faces per frame than produced by the video file used previously, which yields 0.64 faces per frame on average; and (2) because the rate of face production is constant, we do not have to provision our cluster to handle sudden spikes in traffic, allowing us to deploy fewer identification instances than for the video file. None of the conclusions we draw from these experiments is invalidated by these two observations.

Figure 4 shows the effects of accelerating the AI components of *Face Recognition*. Note that because we assume one face per frame, which is significantly higher than the 0.64 faces per frame average produced by our default video file, the average end-to-end latency is somewhat higher at 1× speed than in Section 3.2. At higher speedups, we see a two-fold benefit: first, the latency is very clearly reduced; second, the throughput is commensurately increased.

At 8× speedup, we see a new manifestation of the AI tax, with latency tending toward infinity—the longer the experiment runs, the larger the latency grows. This is an example of an unstable system in queueing theory: faces are entering the system more quickly than they are leaving. This is a major limitation that can severely hamper the prospects of AI acceleration and demands further investigation.

#### 4.4 Network and Storage Bandwidth Limits

With state-of-the-art industry accelerators claiming improvements of up to 15× in inference speedup over CPUs [1], it is critical to understand why the system becomes unbalanced at 8× acceleration. Without this insight, it will be difficult to build systems to accommodate higher acceleration factors.

Intuitively, we suspect the imbalance results from the increased throughput of the system overwhelming one of two resources with limited bandwidth: either network or storage bandwidth. We measure the utilization of both bandwidths to understand the problem at increased acceleration

factors (i.e. greater than 8×). In Figure 5a, the network bandwidth utilization of all container types rises with increasing acceleration factor. Unsurprisingly, producer (ingest/detect containers) network read bandwidth is next to zero, as is consumer (identification containers) write bandwidth. Conversely, producer write and consumer read bandwidths are comparable. But the real network bandwidth hot spot is the brokers—as the point of communication between producers and consumers, they must process all network traffic generated by the producers or read by the consumers. However, even the combined network traffic flowing through the brokers constitutes a small portion of the available bandwidth: at 8× accelerated AI, the read bandwidth is only 6 Gbps, a mere 6% of the available 100 Gbps.

Figure 5b shows the storage bandwidth requirements of the brokers. We omit the data for the producer and consumer containers, as their storage behaviors are not preserved by our emulation technique and are nevertheless expected to be near zero, as they work largely out of memory. The brokers, however, have rather high bandwidth requirements. Even at native (1×) speed, the write bandwidth is 10% of capacity (1.1 GB/s). At 8× acceleration, that rises to over 67%.

With the overhead of the operating system, managing the file system, and coordinating all the small requests to be written to storage, by 8× acceleration, 67% utilization has effectively saturated the available bandwidth. Returning to queueing theory, the inability of storage to write data to storage (and make it available to the consumers) as fast as it is supplied leads to the imbalance and growing latency.

We note that data reads, however, use essentially none of the available bandwidth. This is easily understood: brokers are tasked with ensuring data reliability, so they must write producer data to storage, but the operating system can also cache the data in memory, allowing reads directly from memory and bypassing the storage read path.

Doubtless, fine-tuning the brokers’ parameters could allow them to better utilize the storage bandwidth. An in-depth exploration of the Apache Kafka parameter space is not, however, the purpose of this paper. Regardless of the ability of the brokers to utilize available bandwidth, they will hit a hard limit at the specifications of the hardware devices.

In a setup with a more conservative network bandwidth



(e.g. 10 Gbps), both the storage and the network would quickly become bottlenecks when accelerating compute.

Thus the increased throughput of a moderately accelerated end-to-end system creates a new AI tax that quickly overwhelms the communication substrate, counteracting the gains achieved through hardware acceleration of AI.

#### 4.5 Increase in Waiting Time

Furthermore, whereas the waiting time at  $1\times$  speed constitutes 64.6% of the total latency of a frame (Figure 4), it grows to 66.4% at  $2\times$ , 68.0% at  $4.0\times$ , and 79.1% at  $6\times$ . This trend can be partially understood by Kafka’s automatic batching between brokers and consumers and producers. A message from a producer can be held in the producer for a small amount of time until a larger group of messages has been accumulated to be sent as a batch. Similarly, when a consumer requests available messages from a broker, the broker can withhold messages until there exists some minimum amount of data. Thus, the broker time grows with the decrease in compute time to improve batching. Both batching behaviors are limited by timeouts to ensure that neither producer nor consumer waits excessively long. We have tuned these parameters to find settings that ensure good behavior across a variety of experiments. Nevertheless, the time spent waiting between producers and consumers approaches some lower limit beyond which no amount of tuning can help; in an application that has many more stages than *Face Recognition*, this minimum waiting time could accumulate across stages and prove prohibitively long.

### 5 AI-Centric Data Center Design

As future accelerators emerge, we seek to unlock higher speedups. In Section 4, we saw that in an accelerated AI environment, the AI tax overwhelms the communication mechanism; in particular, the storage medium is quickly saturated at relatively modest emulated compute acceleration speeds. So in Section 5.1 we explore two avenues to overcoming this bottleneck. Implementing these solutions to the tax translates to actual monetary cost, as shown in Section 5.2. We show in Section 5.3 that a purpose-built edge data center can address the tax with increased capacity to handle accelerated compute at lower total cost of ownership (TCO).

#### 5.1 Unlocking Higher Speedups

There are three ways to deal with the limitation in the storage bandwidth: (1) utilize faster storage in the existing brokers, either through a faster storage medium (e.g. Intel Optane [53]) or through multiple drives operating in parallel; (2) create more storage bandwidth by allocating additional brokers; or (3) decrease the size of the face thumbnails, thus demanding less bandwidth. We explore all three methods (Figure 6), first increasing the installed drive count from one to four to provide greater bandwidth to each broker node, then increasing the broker count from three to eight across distinct broker nodes, and finally decreasing the face thumbnails down to one-eighth their original size.

**Increasing the Bandwidth.** The effect of additional storage bandwidth on the existing nodes is captured in Figure 6a. For these experiments, we instantiated additional broker instances on each broker node (one for each drive) to ensure that each drive is given the same access to compute and

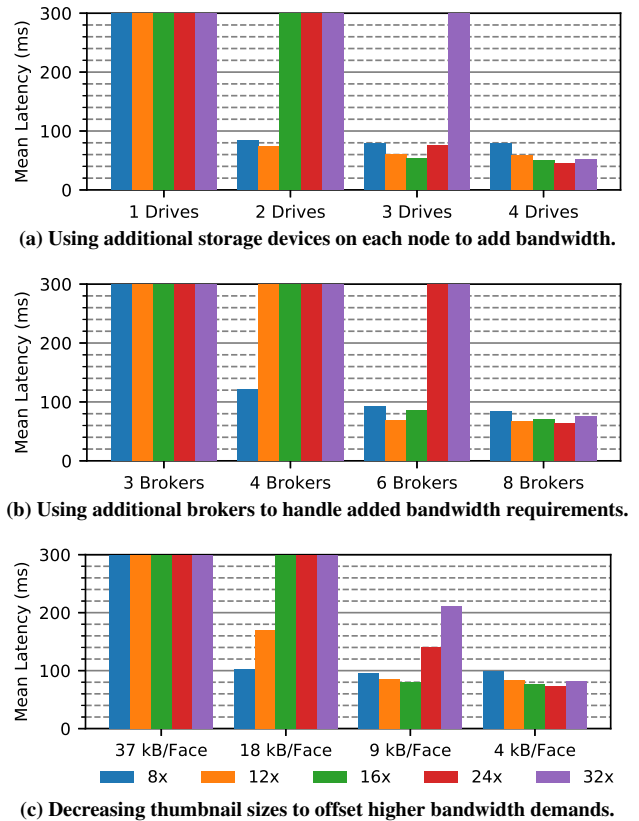


Figure 6: Average frame latency under accelerated AI. The higher bandwidth demands can be accommodated by allocating extra bandwidth or reducing face thumbnail sizes.

memory resources; in practice, only one broker should be instantiated per node to avoid replicating data on the same node. In the figure, we start with  $8\times$  speedup—the speedup that sent latency to infinity in the previous experiments—and increase the emulated speedup to  $32\times$ . With just one NVMe drive, the average end-to-end frame latency is infinite (depicted by the latency bar extending beyond the limits of the chart) at  $8\times$  and all higher speedups. These experiments rely on additional drives being installed only in the brokers—in our case, there are three of them; the remainder of the servers remain unaltered from their original configuration. In increasing the storage bandwidth by going from one drive to two drives, both  $8\times$  and  $12\times$  speedups are “unlocked”—the system gains the ability to support compute of these speeds. With three drives, the system supports up to  $24\times$  speedup, and with four drives,  $32\times$  is unlocked.

**Spreading the Load.** Rather than installing additional drives in each of the brokers, we can instantiate additional brokers in the data center. This spreads out the load on storage to more brokers and hence more drives. Returning each broker to its default storage configuration (one NVMe drive), we repeat our experiments with four, six, and eight brokers.

With three brokers, as we saw before, any acceleration factor at or above  $8\times$  leads to infinite latency. A small 33% increase in the broker count (going from three to four brokers), however, allows the system to handle the  $8\times$  factor,

**Table 2: Homogeneous data center equipment.** In a homogeneous data center similar to (but larger than) our own, there is considerable expense in ensuring that all components are equally equipped. The equipment cost of a 1024-node data center would be around US\$30.9 million.

| Component   | Price (US\$)        | Quantity |
|---|---------------------|----------|
| Dell PowerEdge R740xd (base server)                             | \$28,731            | 1024     |
| Intel Xeon Platinum 8176  | Included            | 2        |
| 32 GB DDR4 SDRAM  | Included            | 12       |
| Intel SSD DC P4510 1 TB (NVMe SSD)                              | \$399               | 1        |
| Mellanox MCX415A (100 GbE adapter)                              | \$660               | 1        |
| Mellanox MSN2700-CS2F<br>(100 GbE switch for fat-tree topology) | \$17,285            | 160      |
| Mellanox MCP1600 (100 GbE cable)                                | \$100               | 3072     |
| <b>Total</b>  | <b>\$33,577,760</b> |          |

while a  $2\times$  broker increase allows for up to  $16\times$  acceleration. At eight brokers, the system can handle a  $32\times$  factor.

We find an important distinction between adding additional drives to existing brokers and adding additional brokers: the latter is more efficient. To achieve the ability for the system to support  $32\times$  accelerated AI compute, we had to increase the number of drives by a factor of four; in contrast, we had to increase the broker count by  $2.7\times$  (going from three brokers to eight) for the same performance achievement. The significantly lower increase in storage bandwidth in the increased-brokers approach indicates that brokers may also benefit from having additional compute capacity, memory bandwidth, or network bandwidth available.

#### Decreasing the Demand.

There exists one additional possibility for reining in the bandwidth demands on the storage: decrease the volume of data that needs to be stored. Rather than spreading the data among additional brokers, the data volume can be reduced by decreasing the average size of face thumbnails. Figure 6c shows the effect of face sizes at one-half, one-quarter, and one-eighth their original size. Similar to increasing the bandwidth in each broker, we see that the smaller face sizes use a smaller portion of the available bandwidth and so increase the maximum supportable speedup, but without instantiating additional brokers or installing additional storage devices.

This solution, however, comes with serious trade-offs. Decreasing face size using compression would require additional compute time, potentially offsetting much or all of the accelerator gains. Decreasing face size by using smaller thumbnails changes the algorithm and can detrimentally impact accuracy. Due to these severe limitations of this approach, we will focus on the previous two solutions.

## 5.2 The Cost of the AI Tax in the Data Center

A typical and simple approach that customers rely on to build an edge data center is to aim for homogeneity across servers (i.e., all of the server components are literally identical across the machines). But in a specialized application domain, such as edge video analytics, this ignores the unique characteristics of the applications and either significantly over-provisions some resources or severely handicaps application performance, leading to suboptimal TCO.

Table 2 shows the basic computing and networking equipment needed to build a homogeneous 1024-node edge data center similar in compute capabilities to our own setup. This design gives each node comparable equipment to that used in our experiments: two 28-core processors, 384 GB of RAM,

100 Gbps interconnect, and a single NVMe drive. The nodes are connected in a three-level fat-tree topology using 32-port Mellanox Ethernet switches. This topology ensures full-speed non-blocking network connectivity to each node.

Using an open source TCO calculator from Coolan [54] to include power (servers, networking equipment, cooling, etc.), rack equipment, cabling costs, etc. and assuming a three-year amortization life, we estimate a yearly cost of US\$10.2 million for server equipment, US\$1.3 million for network equipment, and US\$1.4 million for power, for a total yearly cost of US\$12.9 million.

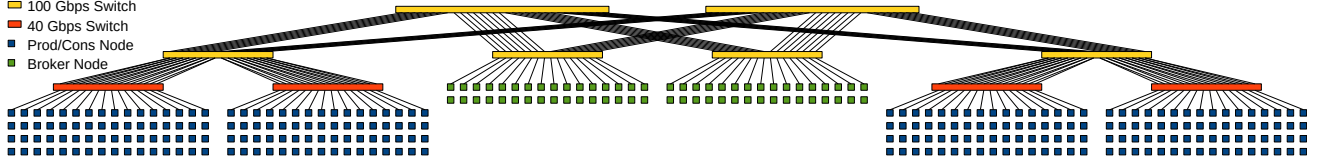
While common wisdom regarding data centers suggests that the majority of the TCO is spent on power (including powering cooling equipment), simple analysis shows this is not necessarily the case. Each of the servers in our hypothetical data center is equipped with a 750 watt power supply, while Mellanox reports that its routers can consume a maximum of 398 watts [55]. This yields a total maximum power consumption of 921 kW. Cooling is estimated to require approximately as much power as the compute resources [56, 57], bringing the total to 1842 kW. Assuming US\$0.10 per kilowatt hour, operating the data center would cost US\$184 per hour or US\$1.61 million per year under maximum load.

To accommodate up to  $32\times$  accelerated compute in AI, we must either install three additional drives in each node (to maintain homogeneity) or designate a large number of the nodes as brokers. Adding the additional NVMe drives costs US\$1.23 million. Instead, we designate 157 of the nodes as brokers, 289 as producers, and 578 as consumers. This maintains the ratio of each node type as in our original *Face Recognition* experiments (15 producer and 30 consumer nodes, though with 8 brokers instead of 3) to enable support for  $32\times$  accelerated AI. Extrapolating from Figure 5a, we estimate each producer and consumer node will consume approximately 4 Gbps of network bandwidth and each broker node 24 Gbps. The broker nodes demand less than 9 Gbps (or 1.1 GB/s) of storage write bandwidth.

## 5.3 AI-Specific Edge Data Center

The homogeneous data center was designed to be generic, capable of executing a variety of application classes; hence, we had to adapt the application to the data center, resulting in hugely over-provisioned network and storage. The producers and consumers constitute over 84% of the data center and use only 4% of the available network capacity and essentially none of the storage bandwidth. The brokers use a respectable 24% of the network capacity and basically all of the storage bandwidth but use very little of the compute capacity. This shows extremely inefficient allocation of limited resources in the data center. But we can do better.

Instead of forcing the application to fit into an existing data center, we propose building a data center that fits the application. We recommend a *purpose-built data center* that specifically targets the broker-specific AI tax (the demand for storage bandwidth). The AI tax, if not accounted for, can translate to non-trivial real-world costs that can drastically affect end users' needs. In contrast, by understanding the AI tax and designing to it, we demonstrate that a moderately-sized edge data center can be purpose-built to better address the AI tax while yielding meaningful cost savings.



**Figure 7: Possible network configuration for an accelerated AI-centric edge data center.** This network configuration is fundamentally a fat-tree built from 100 Gbps switches. However, since no node needs the full 100 Gbps bandwidth, we subdivide it. Mellanox offers splitter cables that split, for example, a 100 Gbps connection into two 50 Gbps or four 25 Gbps connections. Each pair of broker nodes shares a 100 Gbps port. Rather than providing each producer/consumer node with 50 or even 25 Gbps, we further subdivide the network connection speed using 32-port 40 Gbps switches. To provide full bandwidth to the 40 Gbps switches, each 100 Gbps port is split into two 50 Gbps connections, both of which are connected to the same 40 Gbps switch. Thus, the aggregate bandwidth provided to a 40 Gbps switch is 800 Gbps, of which the switch can use 640 Gbps. A 100 Gbps switch can connect two 40 Gbps switches. The slower switches use four-way splitter cables providing each producer/consumer node with 10 Gbps network.

**Table 3: Video analytics-targeted data center equipment.** We use network splitter cables to supply 50 Gbps network to the brokers and, in combination with slower switches, 10 Gbps network to the compute nodes. This offers significant savings on network equipment. Further, we only install NVMe drives in the broker nodes, which are equipped with less compute power than the compute nodes.

| Component   | Price (US\$)        | Quantity |
|---|---------------------|----------|
| Dell PowerEdge R740xd (compute server)                        | \$28,731            | 867      |
| Intel Xeon Platinum 8176                                      | Included            | 2        |
| 32 GB DDR4 SDRAM  | Included            | 12       |
| Mellanox MCX411A (10 GbE adapter)                             | \$180               | 1        |
| Dell PowerEdge R740xd (broker server)                         | \$11,016            | 157      |
| Intel Xeon Bronze 3104  | Included            | 2        |
| 32 GB DDR4 SDRAM  | Included            | 12       |
| Mellanox MCX413A (50 GbE adapter)                             | \$395               | 1        |
| Intel SSD DC P4510 1 TB (NVMe SSD)                            | \$399               | 4        |
| Mellanox MSN2700-CS2F (100 GbE switch)                        | \$17,285            | 28       |
| Mellanox MSN2700-BS2F (40 GbE switch)                         | \$10,635            | 14       |
| Mellanox MFA7A20-C010 (optical splitter 100 GbE to 2x 50 GbE) | \$1,165             | 7        |
| Mellanox MC2609130-003 (copper splitter 40 GbE to 4x 10 GbE)  | \$90                | 217      |
| Mellanox MCP7H00-G002R (copper splitter 100 GbE to 2x 50 GbE) | \$140               | 79       |
| Mellanox MFA1A00-C030 (optical 100 GbE interconnect)          | \$515               | 192      |
| <b>Total</b>  | <b>\$27,878,431</b> |          |

In Table 3 we see the equipment needed for this setup, designed to support up to  $32\times$  AI acceleration. In this scenario, we utilize the same highly parallel servers as in Table 2 for producers and consumers but limit the network bandwidth on these nodes to only 10 Gbps and install only basic storage for operating each server. The broker nodes, in contrast, are built on far less parallel but still impressive CPUs, while enjoying 50 Gbps network connections and four NVMe SSDs.

We illustrate in Figure 7 a simple network solution that could provide the designated bandwidths to each server. At its heart, the network is still a fat-tree built from 100 Gbps Mellanox switches, but, using Mellanox splitter cables and slower 40 Gbps switches, broker nodes are provided with 50 Gbps connections while producer and consumer nodes get 10 Gbps connections. A single edge switch can connect 32 broker nodes or 128 producer/consumer nodes. We can thus build the complete data center using a two-level fat-tree of just 28 100 Gbps switches (12 edge and 16 core); seven edge switches connect to a total of fourteen 40 Gbps switches and five connect to the 157 brokers.

In designing this purpose-built data center, we wanted to

avoid limiting potential advancements or upgrades during the lifetime of the data center. We designed it with double the anticipated requirements for network and storage bandwidth. The brokers were designed to accommodate the  $32\times$  speedup in two separate ways. First, we maintained the higher ratio of brokers to compute nodes, just as we did in the homogeneous design; second, we allocated four times the storage devices and bandwidth to each broker. Either one of these solutions on its own would have been adequate to accommodate the compute speedup. Furthermore, by giving 50 Gbps and 10 Gbps network connections to the broker and compute nodes, respectively, we have allowed them to grow to double their anticipated needs. In combination, we have given the data center the ability to adapt to unanticipated application speedups during its intended lifetime.

Our purpose-built data center incurs an equipment cost of US\$27.9 million with a yearly power cost of US\$1.4 million for a three-year amortized yearly total cost of ownership of US\$10.8 million. This is 16.6% lower than the TCO of the homogeneous data center while being better equipped to handle future accelerated compute.

## 6 Discussion

With the flurry of research and focus on advancements in AI and ML algorithms, it is easy to miss the forest for the trees. One of the most exciting and promising ML benchmarking efforts—MLPerf Inference [58, 59]—offers a suite of ML benchmark tasks, including both vision- and language-based tasks. While its creators took great care to ensure it covers a diverse set of ML tasks, MLPerf Inference is so focused on the ML tasks themselves that it entirely ignores the issue of end-to-end deployment. The reality of any complete AI-based application is that it will operate with pre- and post-processing code, often in the context of a data center, and frequently as part of a multi-stage pipeline. Not only is this context often neglected in current research, it tends to be entirely ignored in accelerator-level research. This leads to three key takeaways.

The first takeaway is that much of the pre- and post-processing code, which is often on the critical path for online request processing, relies on the CPU for its execution. These non AI-stages are often overlooked in research despite the fact that they can dramatically impact the accuracy of the AI models. For instance, choosing a different resizing algorithm like `PIL.Image.BILINEAR`, `cv2.INTER_CUBIC`, and `tf.image.resize_bicubic` for better pre-processing on the CPU can cause the accuracy to vary between 76.16% and

76.52% [60]. Dropping accuracy by 0.4% is a non-negligible change. Therefore, we stress the general-purpose CPU's importance as a crucial workhorse of the AI pipeline.

Second, AI in end-to-end applications often involves multiple queries, not just one inference lookup. However, the vast majority of work tends to focus on optimizing inference for uni-model processing. Optimizing for a cascade of different models stitched together into an AI data flow has yet to be investigated. Like us, Hazelwood et al. [6] also hint at the need to understand the end-to-end AI data flow as that can help lead to new optimizations in minimizing data movement and communication, both within and across nodes.

Finally, we demonstrated using our setup that the streaming infrastructure's underlying communication and storage mechanism will be overwhelmed by the huge throughput of data from accelerated AI processing. Therefore, we recommend expanding beyond CPU and AI compute acceleration and looking for opportunities to accelerate the system as a whole, especially including network and storage. As we demonstrated previously in Section 5, the imbalance between compute and storage requirements of AI and communication nodes presents an opportunity for research on edge data centers to specialize their design, reducing TCO while addressing the unique requirements of AI applications.

In summary, our work has been one of the few, if not only, to look at AI and ML execution in the end-to-end context, evaluating the impact not only of the supporting code but also the infrastructure of an edge data center. We show that the AI tax is significant and identify its various sources. We hope that our findings foster new opportunities for future research. More specifically, we must pursue system-level AI insights in addition to the streamlined focus on AI acceleration because, as we show, the two are not orthogonal.

## 7 Related Work

We build on prior work that enabled and rapidly expanded AI and ML applications. Unlike most of the prior work, however, we explore the implications of accelerating AI computation and how it affects an end-to-end application flow. We present related work in five categories: (1) AI and ML benchmarking, (2) integrating AI and ML, (3) end-to-end application flow studies, (4) exploiting heterogeneity, and (5) edge data centers.

**Benchmarking.** MLPerf is one of the leading resources for benchmarking ML-related compute [58, 59]. It provides flexibility for benchmarking a variety of hardware across a variety of ML kernels, but it entirely ignores the issue of end-to-end application behavior and performance. In our work we demonstrate the central importance of understanding the end-to-end application, showing that each ML kernel can constitute a relatively small portion of the pipeline and that truly optimizing ML performance requires a more holistic view of the system.

**Integrated AI.** In presenting the scale and deployment of ML workloads at Facebook, Hazelwood et al. acknowledged the importance of pre-processing data for training and emphasized its stress on storage, network, and CPU [6]. They acknowledged the potentially high latency of inference for top quality models but did not expose the overhead of pre- and post-processing. Nor did they discuss the resource re-

quirements of streaming inference workloads. We emphasize both of these to show how they can pose a barrier to overall performance improvement from AI acceleration.

Microsoft recognizes the importance of latency in the data center particularly as it applies to deep neural networks [61]. Chung et al. presented Microsoft's Project Brainwave which implements DNNs largely in FPGAs distributed throughout the data center, emphasizing the importance of accelerating increasingly complex DNNs [62]. In contrast, this work emphasizes the importance of the enabling code for AI and assumes that accelerating AI and ML will be successful, instead looking at its ultimate impact on the larger workflow.

**End-to-End Application Flows.** Though not specific to AI workloads, Kanev et al. offered a comprehensive look at the trends of warehouse-scale computing at Google [14]. They quantified data center "taxes"—overheads that come with applications but do not directly contribute to the end result, including compression, communication serialization, and remote procedure calls. We show that the brokers act as a tax, coordinating the activities of a distributed application.

Other work has broken down the end-to-end latency of requests, at various levels of granularity, ranging from evaluation of Internet speed and programming language patterns to operating system scheduling and memory latency [63, 64]. This more closely matches our contribution, though our analysis is restricted to latency within the data center and is focused specifically on the common communication base (Apache Kafka) of open source streaming frameworks in an effort to bring perspective to end-to-end AI application flows.

**Heterogeneous Execution.** Prior works sought to exploit the heterogeneity in a data center, producing benefits in speed, energy consumption, and operating costs [65, 66]. Where these papers sought to capitalize on unintentional heterogeneity (arising from workload co-location and, for example, later upgrades), we extol the benefits of intentionally designing an on-premise data center with heterogeneous servers and network. We thereby add hardware cost savings to the existing benefits of data center heterogeneity.

**Edge Data Centers.** Hewlett Packard Enterprise recently demonstrated that cloud-based computing is often not the most cost-effective solution [7]. Their analysis showed for a well-utilized edge data center, TCO can be drastically lower than comparable capabilities in the cloud. Our work extends that idea, showing how the on-premise data center can be specifically tailored to the needs of AI applications.

## 8 Conclusion

It is easy to get caught up in the excitement of AI and ML; this work has brought context to those advancements, elucidating an AI tax, and serves as a call to action to address limiters of performance in realistic, edge data center deployments of AI applications. Streaming AI applications are only possible with the support of pre- and post-processing code, which is far from trivial in both latency and compute cycles and relies almost exclusively on the CPU for all of the processing. AI applications will likely be composed of multiple inference stages, each with its own characteristics and overheads. And the enabling substrate for managing AI applications in a data center sees hot spots in both network and storage that could soon become bottlenecks if not addressed.

## 9 References

- [1] "Habana Homepage - Habana." <https://habana.ai/>.
- [2] T. Chen, Z. Du, N. Sun, J. Wang, C. Wu, Y. Chen, and O. Temam, "Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning," *ACM Sigplan Notices*, vol. 49, no. 4, pp. 269–284, 2014.
- [3] "Google Cloud Platform Blog: Google supercharges machine learning tasks with TPU custom chip." <https://cloudplatform.googleblog.com/2016/05/Google-supercharges-machine-learning-tasks-with-custom-chip.html>.
- [4] "Intel Unveils the Intel Neural Compute Stick 2 at Intel AI Devcon Beijing for Building Smarter AI Edge Devices." <https://newsroom.intel.com/news/intel-unveils-intel-neural-compute-stick-2/>.
- [5] "Deep Learning and Artificial Intelligence Solutions | NVIDIA." <https://www.nvidia.com/en-us/deep-learning-ai/solutions/>.
- [6] K. Hazelwood, S. Bird, D. Brooks, S. Chintala, U. Diril, D. Dzhulgakov, M. Fawzy, B. Jia, Y. Jia, A. Kalro, *et al.*, "Applied machine learning at facebook: A datacenter infrastructure perspective," in *High Performance Computer Architecture (HPCA), 2018 IEEE International Symposium on*, pp. 620–629, IEEE, 2018.
- [7] H. P. Enterprise, "Hpe on-prem vs. amazon web services (aws)," tech. rep., Hewlett Packard Enterprise Company, 09 2018.
- [8] K. Gyarmathy, "How to Reduce Latency Using Edge Computing." <https://www.vxchnge.com/blog/how-data-center-reduces-latency>.
- [9] M. Hannula, "How Hybrid Cloud Simplifies Data Sovereignty Challenges | CIO." <https://www.cio.com/article/3396631/how-hybrid-cloud-simplifies-data-sovereignty-challenges.html>.
- [10] "The rise of edge data centres - data economy." <https://data-economy.com/the-rise-of-edge-data-centres/>.
- [11] "On-Premise Data Centers: Coming Back or Heading Out?." <https://emconit.com/blog/on-premise-data-centers-coming-back-or-heading-out>.
- [12] "Video analytics market to reach usd 25.4 billion by 2026." <https://www.marketwatch.com/press-release/video-analytics-market-to-reach-usd-254-billion-by-2026-cisco-systems-inc-axis-communications-genetec-inc-2019-09-09>.
- [13] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 815–823, 2015.
- [14] S. Kanev, J. P. Darago, K. Hazelwood, P. Ranganathan, T. Moseley, G.-Y. Wei, and D. Brooks, "Profiling a warehouse-scale computer," in *ACM SIGARCH Computer Architecture News*, vol. 43, pp. 158–169, ACM, 2015.
- [15] "Apache Kafka." <https://kafka.apache.org/>.
- [16] A. Woodie, "Understanding Your Options for Stream Processing Frameworks." <https://www.datanami.com/2019/05/30/understanding-your-options-for-stream-processing-frameworks/>.
- [17] R. Onat, "Apache Storm and Kafka Together: A Real-time Data Refinery." <https://hortonworks.com/blog/storm-kafka-together-real-time-data-refinery/>.
- [18] M. Kleppmann, "Apache Kafka, Samza, and the Unix Philosophy of Distributed Data." <https://www.confluent.io/blog/apache-kafka-samza-and-the-unix-philosophy-of-distributed-data/>.
- [19] R. Metzger, "Kafka + Flink: A Practical, How-To Guide." <https://www.ververica.com/blog/kafka-flink-a-practical-how-to>.
- [20] DataTorrent, "End-to-end 'Exactly-Once' With Apache Apex." [https://cdn.rawgit.com/dtpublic/website/b0c73294/blogs/End-to-end%20Exactly-Once\\_%20with%20Apache%20Apex%20-%20DataTorrent.htm](https://cdn.rawgit.com/dtpublic/website/b0c73294/blogs/End-to-end%20Exactly-Once_%20with%20Apache%20Apex%20-%20DataTorrent.htm).
- [21] F. Yang, "Building a Streaming Analytics Stack with Apache Kafka and Druid." <https://www.confluent.io/blog/building-a-streaming-analytics-stack-with-apache-kafka-and-druid/>.
- [22] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, *et al.*, "TensorFlow: Large-scale machine learning on heterogeneous distributed systems," *arXiv preprint arXiv:1603.04467*, 2016.
- [23] C. Regazzoni, A. Cavallaro, Y. Wu, J. Konrad, and A. Hampapur, "Video analytics for surveillance: Theory and practice [from the guest editors]," *IEEE Signal Processing Magazine*, vol. 27, no. 5, pp. 16–17, 2010.
- [24] "Amazon.com: : Amazon Go." <https://www.amazon.com/b?node=16008589011>.
- [25] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multi-task cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.
- [26] "GitHub - davidsandberg/facenet: Face recognition using Tensorflow." <https://github.com/davidsandberg/facenet>.
- [27] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *AAAI*, vol. 4, p. 12, 2017.
- [28] "Docker - Build, Ship, and Run Any App, Anywhere." <https://www.docker.com/>.
- [29] "Apache Storm." <http://storm.apache.org/>.
- [30] "Samza." <http://samza.apache.org/>.
- [31] "Apache Flink: Stateful Computations over Data Streams." <https://flink.apache.org/>.
- [32] "Apache Apex." <http://apex.apache.org/>.
- [33] "Druid | Interactive Analytics at Scale." <https://druid.apache.org/>.
- [34] "Apache Kafka." <https://kafka.apache.org/documentation/streams/>.
- [35] K. Birman and T. Joseph, *Exploiting virtual synchrony in distributed systems*, vol. 21. ACM, 1987.
- [36] C. E. Leiserson, "Fat-trees: universal networks for hardware-efficient supercomputing," *IEEE transactions on Computers*, vol. 100, no. 10, pp. 892–901, 1985.
- [37] "Intel® Xeon® Platinum 8176 Processor (38.5M Cache, 2.10GHz) Product Specifications." <https://ark.intel.com/content/www/us/en/ark/products/120508/intel-xeon-platinum-8176-processor-38-5m-cache-2-10-ghz.html>.
- [38] "Intel® SSD DC P4510 Series (1.0TB, 2.5in PCIe 3.1 x4, 3D2, TLC) Product Specifications." <https://ark.intel.com/content/www/us/en/ark/products/122573/intel-ssd-dc-p4510-series-1-0tb-2-5in-pcie-3-1-x4-3d2-tlc.html>.
- [39] "Production-Grade Container Orchestration - Kubernetes." <https://kubernetes.io/>.
- [40] K. Pingali, "A Case for Case Studies." <https://www.sigarch.org/a-case-for-case-studies/>, 2019.
- [41] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, pp. 91–99, 2015.
- [42] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, *et al.*, "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.
- [43] "Open Source Search & Analytics · Elasticsearch | Elastic." <https://www.elastic.co/>.
- [44] "Logstash: Collect, Parse, Transform Logs | Elastic." <https://www.elastic.co/products/logstash>.
- [45] A. Lindoso and L. Entrena, "Hardware architectures for image processing acceleration," in *Image Processing*, IntechOpen, 2009.
- [46] U. Gupta, X. Wang, M. Naumov, C.-J. Wu, B. Reagen, D. Brooks,



- B. Cotel, K. Hazelwood, B. Jia, H.-H. S. Lee, *et al.*, “The architectural implications of facebook’s dnn-based personalized recommendation,” *arXiv preprint arXiv:1906.03109*, 2019.
- [47] D. Mahajan, J. Park, E. Amaro, H. Sharma, A. Yazdanbakhsh, J. K. Kim, and H. Esmailzadeh, “Tabla: A unified template-based framework for accelerating statistical machine learning,” in *2016 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 14–26, IEEE, 2016.
- [48] Xilinx, “Accelerating dnn with xilinx alveo accelerator cards,” tech. rep., Xilinx, Inc., 10 2018.
- [49] “NVIDIA Deep Learning Accelerator.” <http://nvidia.org/>.
- [50] “Coral.” <https://coral.withgoogle.com/>.
- [51] “Inference - Habana.” <https://habana.ai/inference/>.
- [52] B. Wheeler, “Data centers accelerate ai processing,” tech. rep., The Linley Group, 12 2018.
- [53] “Intel® Optane™ Technology.” <https://www.intel.com/content/www/us/en/architecture-and-technology/intel-optane-technology.html>.
- [54] “Slash data-center costs and downtime by using Coolan’s TCO Model - TechRepublic.” <https://www.techrepublic.com/article/slash-data-center-costs-and-downtime-by-using-coolans-tco-model/>.
- [55] “Specifications - SN2000 Series - Mellanox Docs.” <https://docs.mellanox.com/display/sn2000pub/Specifications>.
- [56] “HPE Power Advisor.” <https://paonline56.itcs.hpe.com/?Page=Index#>.
- [57] “Data Center Cooling Costs | Dataspan.” <https://www.dataspan.com/blog/data-center-cooling-costs/>.
- [58] V. Janapa Reddi, C. Cheng, D. Kanter, P. Mattson, G. Schmuelling, C.-J. Wu, B. Anderson, M. Breughe, M. Charlebois, W. Chou, R. Chukka, C. Coleman, S. Davis, P. Deng, G. Diamos, J. Duke, D. Fick, J. S. Gardner, I. Hubara, S. Idgunji, T. B. Jablin, J. Jiao, T. St. John, P. Kanwar, D. Lee, J. Liao, A. Likhomotov, F. Massa, P. Meng, P. Micikevicius, C. Osborne, G. Pekhimenko, A. T. R. Rajan, D. Sequeira, A. Sirasao, F. Sun, H. Tang, M. Thomson, F. Wei, E. Wu, L. Xu, K. Yamada, B. Yu, G. Yuan, A. Zhong, P. Zhang, and Y. Zhou, “Mlperf inference benchmark,” 2019.
- [59] “MLPerf.” <https://mlperf.org/inference-overview/>.
- [60] B. Felter, “[green] change preprocess module from PIL to opencv by ClarkChin08 · Pull Request #273 · mlperf/inference · GitHub.” <https://github.com/mlperf/inference/pull/273>.
- [61] “Project brainwave - microsoft research.” <https://www.microsoft.com/en-us/research/project/project-brainwave/>.
- [62] E. Chung, J. Fowers, K. Ovtcharov, M. Papamichael, A. Caulfield, T. Massengill, M. Liu, D. Lo, S. Alkalay, M. Haselman, *et al.*, “Serving dnn in real time at datacenter scale with project brainwave,” *IEEE Micro*, vol. 38, no. 2, pp. 8–20, 2018.
- [63] M. Chow, D. Meisner, J. Flinn, D. Peek, and T. F. Wenisch, “The mystery machine: End-to-end performance analysis of large-scale internet services,” in *11th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 14)*, pp. 217–231, 2014.
- [64] J. Li, N. K. Sharma, D. R. Ports, and S. D. Gribble, “Tales of the tail: Hardware, os, and application-level sources of tail latency,” in *Proceedings of the ACM Symposium on Cloud Computing*, pp. 1–14, ACM, 2014.
- [65] J. Mars, L. Tang, and R. Hundt, “Heterogeneity in “homogeneous” warehouse-scale computers: A performance opportunity,” *IEEE Computer Architecture Letters*, vol. 10, no. 2, pp. 29–32, 2011.
- [66] M. E. Haque, Y. He, S. Elnikety, T. D. Nguyen, R. Bianchini, and K. S. McKinley, “Exploiting heterogeneity for tail latency and energy efficiency,” in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 625–638, ACM, 2017.