

## DEPARTMENT: VIEW FROM THE CLOUD

# Sustainable and Trustworthy Edge Machine Learning

Ivona Brandić , Vienna University of Technology, 1040 Vienna, Austria

*Nowadays, our world is driven by complex, large scale, yet tactile information systems requiring various degrees of trustworthiness. Trustworthiness of the systems always comes with costs. The traditional and rather costly way to understand the behavior of large scale systems is to develop powerful mathematical abstractions that allow us to condense these behaviors and to reason about them at a very abstract level. In our FWF funded project Rucon, we introduce an orthogonal, data driven, and probabilistic concept to model and reason uncertainty of the systems. In Rucon, deliberated system failures are tolerated due to the benefits of the costs and sustainability. Rucon's approach targets large scale near real-time systems like live video analytics, streaming, vehicular applications, and smart city information systems.*

The ongoing digital transformation is disruptively changing all aspects of our lives. The sectors immediately affected and revolutionized by the IoT are healthcare (smart medical devices), manufacturing (smart factories), energy (smart power grids) as well as urban development and transportation (smart buildings, cities and vehicles). One of the direct impacts of the increased digitalization is the nearly exponential increase of energy demand to process and store data. Already, data centers represent an estimated 1% of global electricity demand. One of the most worrying models predicts that electricity use by ICT could exceed 20% of the global total by the time a child born today reaches her teens.<sup>9</sup>

To meet the demands of the ongoing digitization efforts, a new generation of information systems is emerging with latencies less than 100 ms or even less than 10 ms what is called nowadays “tactile internet,” addressing upcoming data driven business applications like virtual reality, telemedicine, smart cities, or self-driving cars. The trend in transforming traditional backend applications to “tactile internet” applications is also affecting the High Performance Computing

(HPC) area. In the area of HPC, we have the concept of “Extreme Data” that follows the “Big Data” problem, by addressing massive amounts of information that must be processed and analyzed in near-real time through the utilization of Exascale systems. In the past decade, Cloud Data Centers and Supercomputers have been envisioned as the essential computing architectures for enabling the next generation of extreme data applications. However, in recent years we experienced the rise of near-real time extreme data systems. These applications process complex data intensive workflows with strict latency requirements. An example of an extreme data application is the early earthquake alert based on the analysis of thousands of sensors from smart phones.<sup>7</sup> In both cases (commercial and HPC computing), we evidence the paradigm shift not only to “tactile internet” but also in the type of emerging applications, where traditional simulation and optimization applications are replaced or enhanced with data intensive machine learning (ML) applications.

When considering such huge, complex, and geographically distributed systems, it is not intuitively clear where the most common performance bottlenecks are, or which parts of the systems are the most inefficient in terms of the energy consumption. In case of the tactile internet, both the hardware and the application can contribute equally to the inefficiency and high failure rates.

## DATA INTENSIVE APPLICATIONS

Data intensive applications in the context of tactile internet differ fundamentally from traditional ML applications trained and executed in centralized and highly controlled environments like Cloud Data Centers.<sup>14</sup> Many geographically distributed ML applications, e.g., Apple Siri are entirely based on cloud computing. Such applications do not function if the network is unavailable. Also many of the existing intelligent applications generally adopt centralized data management, where users upload their data to a central cloud based data center. However, with the ever increasing volumes of data, which has been generated and collected by billions of mobile users and IoT devices it is estimated that Internet traffic is reaching 235.7 Exabytes per month in 2021, up from 73.1 Exabytes per month in 2016.<sup>2</sup>

## HYPERHETEROGENEOUS HARDWARE

Hardware architectures and networks utilized for IoT systems and tactile internet differ radically from all other well known systems. Since IoT devices (e.g., sensors) are rather tiny and not capable of running complex computation, more powerful nodes in the vicinity of IoT devices are necessary to process and store data and ensure low latency; this is called “the Edge.” The concept, when the local version of the ML model is deployed at the Edge, is called *Edge Machine Learning (EML)*.<sup>1,10,11</sup> An important application area for EML is (near) real-time object detection, as it is currently often impossible to run the video inference without GPU on board.<sup>5</sup> Another application area is to enhance 5G with computational facilities.<sup>6</sup> The resource landscape is becoming more and more heterogeneous with Edge nodes that can significantly vary in their shape, size, and computational power, resulting in the so-called hyperheterogeneity. Edge nodes might range from the so-called  $\mu$ -Data Centers<sup>8</sup> consisting of several servers to simple Raspberry Pis. Network connections might vary as well ranging from wifi to LTE or 4G.

## CHALLENGES IN TERMS OF SUSTAINABILITY AND TRUSTWORTHINESS

Hyper heterogeneity and geographical distribution of Edge systems make it difficult to manage the competing priorities like sustainability and trustworthiness. Even worse, a high degree of geographical distribution very often results in intermittent connectivity that prevents us from utilizing well-known sustainability and trustworthiness concepts from Cloud Computing or other types of distributed systems. A typical concept

to achieve sustainable systems in Clouds is to shut down virtual machines (VMs) in case of low workload. Edge systems very often rely on event driven microservices that do not allow management of resources at the granularity level of VMs. A similar example for the lack of trustworthiness is the requested availability of Edge systems. In Cloud systems, we can achieve high availability by using sophisticated backend scheduling and load balancing algorithms. Sophisticated scheduling and load balancing are difficult at the Edge due to resource scarcity.

In our FWF Rucon (Runtime Control on Multi-Clouds) project,<sup>a</sup> we developed in the last five years the fundamentals for sustainable and trustworthy Edge Machine Learning systems with two major goals:

- 1) *ML for the Edge*: The first goal is to develop ML based methods for the sustainable and trustworthy operation of Edge nodes, regardless of the applications being executed at the particular Edge node.
- 2) *ML at the Edge*: The second goal of the Rucon project is to develop methods for the sustainable and trustworthy execution of geographically distributed ML applications (e.g., streaming apps).

In order to make decisions on the huge amounts of data in a relatively short time frame, the whole Rucon architecture was developed in a probabilistic manner capable of dealing with hyper heterogeneity, geographical distribution, intermittent connectivity, and low availability of the nodes. The centerpiece of the Rucon architecture are as follows:

- › The novel method for the fault tolerance and trustworthiness based on the Dynamic Bayesian Networks (addressing ML for the Edge).
- › A sustainable model management approach based on the Reinforcement Learning (RL) for geographically distributed ML (addressing ML at the Edge).
- › A sustainable and trustworthy method for data quality management for failure prone IoT devices (addressing both ML at the Edge and ML for the Edge).

## TRUSTWORTHINESS IN RUCON: FAULT TOLERANCE

Edge computing is prone to failures as it trades reliability against other QoS properties such as low latency and geographical prevalence. Failures on the Edge happen much more often than in other large scale systems

<sup>a</sup><http://rucon.ec.tuwien.ac.at>



**FIGURE 1.** (a) Object detection. (b) Traffic light with the camera and Raspberry Pi. (c) App with the alert about the object in the dead corner.

due to geographical dispersion, ad hoc deployment, and rudimentary support systems (e.g., lack of diesel generators to compensate for power outages). Software services that run on Edge infrastructures must rely on failure resilience techniques for uninterrupted delivery. Due to the lack of other support systems, this has to happen at the software layer. Edge nodes are usually deployed in urban areas with space restriction and using low cost devices like well known Raspberry Pis. Figure 1(a)–(c) depicts such a real-life system installed to collect traces and data samples for our experimental smart traffic light system.<sup>b</sup>

We developed an app that visualizes objects on the smartphone that appear in vehicles' dead corners, thus preventing severe accidents. To detect objects [Figure 1(a)], we used cameras inside a traffic light [Figure 1(b)], together with a Raspberry Pi equipped with convolutional networks for object detection. Once an object is detected in the dead corner, the message is broadcast to all vehicles of interest while the detected objects are visualized on the app as shown in Figure 1(c).<sup>13</sup> As can be seen on the pictures, we used low cost devices suitable for mass rollout in smart cities. However, these devices can easily fail and require sophisticated failure tolerance mechanisms. A well-known approach to counteract low availability is the utilization of geographically distributed replicas for the deployed services. In case of the failure of a service, the workload is redistributed to the standby replica. Standby replicas, however, should not fail concurrently.

In *Rucon*, we developed a novel fault-tolerance mechanism for the redundant service deployment that minimizes the cost (e.g., in terms of the number

of redundant services) while preventing joint failures of the replicas. Spatiotemporal dependencies of failures appear very frequently in Edge systems. Reasons might be a network failure affecting multiple servers in the same physical/virtual network or a power outage affecting multiple servers in the same grid or multiple servers deployed in hostile locations failing due to environmental interference. Neglected spatiotemporal dependencies can lead to the so-called cascading failures, and in general to catastrophic effects for the overall reliability of systems. In *Rucon*, we detect spatiotemporal failure dependencies among Edge servers to improve the failure resilience of services with minimum possible redundancy by applying the dynamic Bayesian networks (DBNs). In this architecture, dependence learning occurs in a resource-rich environment such as the cloud based on the received past failure traces of the system. Trained DBNs are then used to perform the inference about the joint failure probability of the random servers.<sup>3</sup>

In our approach, we learn the spatiotemporal dependencies between Edge server failures and combine them with the topological information to incorporate link failures. Eventually, we infer the probability that a certain set of servers fails or disconnects concurrently during service runtime. Our experimental results show that after eliminating the noise and by analyzing randomly large scale failure traces of Edge datasets of various applications (e.g., Skype supernodes), there is a significant amount of spatiotemporal failures. We developed multiple dependence- and topology-aware deployment algorithms that minimize either failure probability or redundancy cost. Experimental results show that we can reduce the service downtime by several orders of magnitude compared to the baseline while preserving the requested latency. The utilization of the deployment

<sup>b</sup><http://intrasafed.ec.tuwien.ac.at>

algorithms that consider the joint failure probability can further decrease the redundancy loss up to 50% compared to the baseline. We consider our spatiotemporal failure dependence approach as the first step toward trustworthy edge machine learning.

### SUSTAINABILITY IN RUCON: MODEL DISTRIBUTION FOR EDGE MACHINE LEARNING

In a distributed setting, ML models are usually trained in a large scale data center. Afterward (a reduced) version of the model is distributed to the Edge to perform inference in the vicinity of the end users and thus achieve low latency. When deploying ML models over geographically distributed Edge nodes several problems arise, in particular in nonstationary environments when the data distribution changes. Due to environmental changes models that have been learned, trained, and distributed to Edge nodes might become inaccurate and in the worst case no longer valid. In traditional data centers, nonstationarity is solved using so called online learning, where models are trained in batches as new data arrives. Applying online learning in a geographically distributed setting bears several problems in terms of sustainability, where distributed ML models can be independently trained and periodically synchronized through a centralized parameter server. Too frequent updates would result in a heavy message exchange and could lead to bandwidth problems and eventually bad sustainability of the whole system. Less frequent model updates can result in poor performance of the ML models at the Edge.

This dilemma might occur in many geographically distributed streaming applications, which are very common in the area of the “tactile internet.” One such example is the e-vehicle application managing recharging intervals of electrical cars. Changing environmental conditions (wind, sun, water, etc.,) cause volatile availability of the electricity, which has to be matched with user requests for short queuing intervals at the charging stations. In this setup, cars communicate with roadside units (RSU) to be updated about the current situation at the charging stations, i.e., to receive the freshest version of the ML model. On the other hand, RSUs are used to collect data from the cars to accurately predict the needs of car fleets on the road.

With our staleness control mechanism proposed in *Rucon*, we tackle the concept drift issues in Edge data analytics to minimize its accuracy loss of the distributed ML without losing its timeliness benefits. We propose an efficient model synchronization mechanism for distributed and stateful data analytics. Our RL-based algorithm learns over time the connectivity patterns of the cars and the most suitable intervals for the distribution of newly

collected data in the form of model updates at the RSU. The data are distributed from the car to the RSU on one hand, and also from the parameter server to the RSU in the form of the updated models. Since we use online RL, the algorithm has a low computational overhead, automatically adapts to changes, and does not require additional data monitoring contributing to the sustainability of Edge Machine Learning. Our thorough evaluation shows that we are able to save up to 90% of the updates while having the same quality of the model compared to the fully synchronous oracle approach.<sup>4</sup> Reduced number of model updates directly increases the energy efficiency of the geographically distributed ML applications.

### SUSTAINABILITY AND TRUSTWORTHINESS IN RUCON: DATA QUALITY MANAGEMENT

Edge computing is usually utilized to collect and process data from the sensors and other IoT devices that have a very high failure rate. Missing or invalid data may appear very often on the IoT systems due to monitoring system failures, data packet loss, or sensor aging. Consequently, near-real-time decisions are often based on limited and incomplete data. Low data quality might significantly impact accuracy of the decision-making processes on a large scale, e.g., in large scale cloud data centers that use collected, aggregated, and processed IoT data.

*Rucon*’s approach to counteract the low data quality is a generic mechanism for recovery of multiple gaps in incomplete datasets, using multiple recovery techniques. To ensure outliers removal, detection, and forecasting of each gap, we use different techniques addressing different dataset characteristics. The autoregressive integrated moving average (ARIMA) method can be used, if data contain stationary characteristics, such as trend stationarity. The Exponential Smoothing method (ETS) can be used for short-term seasonal series or with multiple complex seasonality. Another feature of *Rucon* is a sustainable Edge data management that achieves a tradeoff between the amount of data stored at the Edge and high accuracy for predictive analytics. Our data quality approach facilitates adaptive storage management mechanisms for reducing the amount of data stored at the Edge, keeping only the data necessary for predictive analytics.<sup>12</sup> We utilize data clustering techniques where we detect stable accuracy clusters. Those clusters can be used as a border between relevant and irrelevant data for the accurate near-real time analytics. Once identified, irrelevant data can be released and thus data storage is optimized on the resource scarce Edge nodes. Data quality management addresses both, sustainability of the Edge storage but also trustworthiness of processed data.



## DISCUSSION AND OUTLOOK

In *Rucon*, we have developed the first fundamental approaches for achieving sustainable and trustworthy Edge Machine Learning focusing on the current challenges like hyper heterogeneity and high failure rate of IoT sensors. However, the demand for the sustainable and trustworthy Edge systems will significantly increase in the future as discussed next.

**Arbitrary resources.** Nowadays, we usually have dedicated Edge nodes, for example, installed at road side units or in combination with 5G antennas. Usually, 5G antennas are equipped with additional servers to process the workload on demand in the vicinity of the end users. For moving objects like drones and scooters, it is even harder to facilitate efficient resource usage as they have intermittent connectivity and very frequent handovers between the Edge nodes. Building stationary Edge nodes, where drones fly every now and then, is highly inefficient. The future challenge is to develop sustainable and trustworthy Edge systems even in case of arbitrary and/or opportunistic computing. In both cases, the idea is to incentivize people to share their resources if there is a high demand for them. The computation could be dynamically offloaded and dynamically migrated on already existing but idle systems (e.g., idle laptop, idle server) in a secure way, if the required middleware for the management and charging of such systems is installed. This paradigm will only succeed, if the resources are trustworthy. On the other hand, arbitrary and/or opportunistic computing is the basic concept of the sharing economy and can create many other benefits in terms of sustainability in the long run. 🌱

## ACKNOWLEDGMENTS

This work was supported in part by the Austrian Science Fund (FWF Y 904 START-Programm 201) and in part by the City of Vienna (5G Use Case Challenge InTraSafEd 5G).

## REFERENCES

1. Okanojara, Daisuke, *et al.*, "Machine learning with model filtering and model mixing for edge devices in a heterogeneous environment," U.S. Patent No. 10 387 794, Google Patent, 0217387A1/en, 2016. [Online]. Available: <https://patents.google.com/patent/US>
2. CISCO, *VNI Complete Forecast Highlights*, 2016. [Online]. Available: [https://www.cisco.com/c/dam/m/en\\_us/solutions/service-provider/vni-forecasthighlights/pdf/Global\\_2021\\_Forecast\\_Highlights.pdf](https://www.cisco.com/c/dam/m/en_us/solutions/service-provider/vni-forecasthighlights/pdf/Global_2021_Forecast_Highlights.pdf)

3. A. Aral and I. Brandic, "Learning spatiotemporal failure dependencies for resilient edge computing services," *IEEE Trans. Parallel Distrib. Syst.*, vol. 32, no. 7, pp. 1578–1590, Jul. 2021.
4. A. Aral, M. Erol-Kantarci, and I. Brandic, "Staleness control for edge data analytics," *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 4, no. 2, pp. 38:1–38:24, 2020.
5. L. Cavigelli, P. Degen, and L. Benini, "CBInfer: Change-based inference for convolutional neural networks on video data," in *Proc. 11th Int. Conf. Distrib. Smart Cameras*, Stanford, CA, USA, 2017, pp. 1–8.
6. Y. C. Hu *et al.*, "Mobile edge computing: A key technology towards 5G," *ETSI White Paper*, vol. 11, no. 11, pp. 1–16, 2015.
7. K. Fauvel *et al.*, "A distributed multi-sensor machine learning approach to earthquake early warning," in *Proc. 34th AAAI Conf. Artif. Intell./32nd Innov. Appl. Artif. Intell. Conf./10th AAAI Symp. Educ. Adv. Artif. Intell.*, 2020, pp. 403–411.
8. A. G. Greenberg, J. R. Hamilton, D. A. Maltz, and P. Patel, "The cost of a cloud: Research problems in data center networks," *Comput. Commun. Rev.*, vol. 39, no. 1, pp. 68–73, 2009.
9. N. Jones, "How to stop data centres from gobbling up the world's electricity," *Nature*, vol. 561, no. 7722, pp. 163–167, Sep. 2018. [Online]. Available: <https://www.nature.com/articles/d41586-018-06610-y>
10. H. Li, K. Ota, and M. Dong, "Learning IoT in edge: Deep learning for the internet of things with edge computing," *IEEE Netw.*, vol. 32, no. 1, pp. 96–101, Jan./Feb. 2018.
11. D. Liu, G. Zhu, J. Zhang, and K. Huang, "Data-importance aware user scheduling for communication-efficient edge machine learning," *IEEE Trans. Cogn. Commun. Netw.*, vol. 7, no. 1, pp. 265–278, Mar. 2021.
12. I. Lujic, V. De Maio, and I. Brandic, "Resilient edge data management framework," *IEEE Trans. Serv. Comput.*, vol. 13, no. 4, pp. 663–674, Jul./Aug. 2020.
13. I. Lujic, V. De Maio, K. Pollhammer, I. Bodrozic, J. Lasic, and I. Brandic, "Increasing traffic safety with real-time edge analytics and 5G," in *Proc. 4th Int. Workshop Edge Syst., Analytics Network.*, 2021, pp. 19–24.
14. M. Zaharia *et al.*, "Apache spark: A unified engine for big data processing," *Commun. ACM*, vol. 59, no. 11, pp. 56–65, 2016.

**IVONA BRANDIC** is a Full Professor for High Performance Computing Systems at the Vienna University of Technology. In 2015 she was awarded the FWF START prize, the highest Austrian award for early career researchers. She received the Ph.D. degree in 2007 and the Distinguished Young Scientist Award in 2011, both from the Vienna University of Technology. Her main research interests are runtime management of large scale distributed systems, Cloud Computing, energy efficiency, QoS and autonomic computing. Contact her at [ivona.brandic@tuwien.ac.at](mailto:ivona.brandic@tuwien.ac.at).