Distributed Machine Learning in Edge Computing

Subjects: Computer Science, Artificial Intelligence

Contributor: Carlos Poncinelli Filho , Elias Marques Junior , Victor Chang , Leonardo dos Santos , Flavia Bernardini , Paulo F. Pires , Luiz Satoru Ochi , Flavia Coimbra Delicato

Distributed edge intelligence is a disruptive research area that enables the execution of machine learning and deep learning (ML/DL) algorithms close to where data are generated. Since edge devices are more limited and heterogeneous than typical cloud devices, many hindrances have to be overcome to fully extract the potential benefits of such an approach (such as data-in-motion analytics).

machine learning artificial intelligence distributed edge intelligence fog intelligence

Internet of Things

1. Introduction

Nowadays, with the rise of the Internet of Things (IoT), a large number of smart applications are being built, taking advantage of connecting several types of devices to the internet. These applications will generate a massive amount of data that need to be processed promptly to generate valuable and actionable information. Edge intelligence (EI) refers to the ability to bring about the execution of machine learning tasks from the remote cloud closer to the IoT/Edge devices, either partially or entirely. Examples of edge devices are smartphones, access points, gateways, smart routers and switches, new generation base stations, and micro data centers.

Some edge devices have considerable computing capabilities (although always much smaller than cloud processing centers), but most are characterized by very limited capabilities. Currently, with the increasing development in the area of MEMS (Micro–Electro–Mechanical Systems) devices, there is a tendency to carry out part of the processing within the data producing devices themselves (sensors) ^{[1][2][3][4]}. There are certainly several challenges involved in performing processing on resource-limited devices, including the need to adapt complex algorithms and divide the processing among several nodes.

Therefore, in Edge Intelligence, it is essential to promote collaboration between devices to compensate for their lower computing capacity. Some synonyms of this concept found in the literature are: distributed learning, edge/fog learning, distributed intelligence, edge/fog intelligence and mobile intelligence ^{[5][6][7]}.

The leverage of edge intelligence reduces some drawbacks of running ML tasks entirely in the cloud, such as:

- High latency ^[8]: offloading intelligence tasks to the edge enables achievement of faster inference, decreasing the inherent delay in data transmission through the network backbone;
- Security and privacy issues ^{[9][10]}: it is possible to train and infer on sensitive data fully at the edge, preventing their risky propagation throughout the network, where they are susceptible to attacks. Moreover, edge intelligence can derive non-sensitive information that could then be submitted to the cloud without further processing;
- The need for continuous internet connection: in locations where connectivity is poor or intermittent, the ML/DL could still be carried out;
- Bandwidth degradation: edge computing can perform part of processing tasks on raw data and transmit the produced data to the cloud (filtered/aggregated/pre-processed), thus saving network bandwidth. Transmitting large amounts of data to the cloud burdens the network and impacts the overall Quality of Service (QoS) ^[11];
- Power waste ^[12]: unnecessary raw data being transmitted through the internet demands power, decreasing energy efficiency on a large scale.

The steps for data processing in ML vary according to the specific technique in use, but generally occur in a welldefined life cycle, which can be represented by a workflow. Model building is at the heart of any ML technique, but the complete life cycle of a learning process involves a series of steps, from data acquisition and preparation to model deployment into a production environment. When adopting the Edge intelligence paradigm, it is necessary to carefully analyze which steps in the ML life cycle can be successfully executed at the edge of the network. Typical steps that have been investigated for execution at the edge are data collection, pre-processing, training and inference.

2. Related Work

Some surveys have been published that address the edge intelligence subject recently. However, they adopt different perspectives from the one adopted in this SLR. Al-Rakhami et al. ^[13] propose and analyze a framework based on the distributed edge/cloud paradigm using docker technology which provides a very lightweight and effective virtualization solution. This solution can be utilized to manage, deploy and distribute applications onto clusters (e.g., small board devices such as Raspberry PI). It is able to provide an advantageous combination of various benefits and lower costs of data processing performed at the edge instead of central servers. However, the authors base their proposal on experiments to support the proposal of a new framework. The research does not mention any of the nine groups of techniques the researchers present in the work.

Wang et al. ^[14] survey is centered on the connection between Deep Learning and the edge, either to apply DL in optimizing the edge or to use the edge to run DL algorithms. The study is divided into five fronts: DL applications on edge; DL inference in edge; edge computing for DL; DL training at the edge; DL for optimizing the edge. The paper discusses hardware and virtualization aspects. Concerning the (groups of) techniques and strategies, it is more restricted to Federated Learning and the optimization of the edge with DL.

Xu et al. ^[10] approach edge intelligence under the perspectives of edge caching, edge training, edge inference, and edge offloading in a very comprehensive way. The researchers discuss all these aspects in the work but explore additional techniques, and strategies related to pre-processing, federated learning, and scheduling. One intersection of this paper with the researchers' research is the overlap of three groups of techniques the researchers present (Federated Learning, Edge Pre-processing and Scheduling). However, the researchers deepened the discussion into more groups of techniques.

The work presented by Zhou et al. ^[15] covers artificial intelligence to edge AI, showing a generalized representation of application architecture used in the lifecycle management of ML. In the edge layer: sensors/actuators; edge analytics; logging and monitoring. In the fog layer: visualization; live streaming engines; batch processing; data ingestion; storage and ML model development platforms and libraries. The researchers' research approaches several more domains in which edge intelligence is used, which are not present in this survey. Compared to these other surveys, the researchers analyze the literature more comprehensively, including a discussion on application domains of edge intelligence and their correlation with identified techniques.

Verbraeken et al. ^[16] provide an extensive overview of the current state-of-the-art in terms of outlining the challenges and opportunities of distributed machine learning over conventional machine learning, discussing the techniques used for distributed machine learning. The paper follows the same line of research of Wang et al. ^[14], with a focus on machine learning applied to the distributed environment. To this end, it makes inroads into the various types of algorithms to solve problems using ML.

Table 1 shows the comparison between the researchers' work and the other surveys mentioned in this section. In summary, the main gaps of the analyzed works are focused on aspects such as "Techniques and Strategies" on the edge. The table also shows the aspects of "Challenges" and "Different Application Domains", where edge intelligence can be used.

		Scope	
Paper	Challenges	Group	Different Application
	Challenges	of Techniques	Domains
Al-Rakhami et al. ^[13]	0/6	2/8	1/6
Wang et al. ^{[<u>14]</u>}	1/6	4/8	4/6
Verbraeken et al. ^[16]	1/6	0/8	0/6

Table 1. Comparison of existing surveys.

on

https://www.semi.org/en/blogs/technology-trends/ai-and-mems-sensors (accessed on 1 March 2022).

2.3 ioAnswering the Ros. Reservoir computing with a single delay-coupled non-linear

mechanical oscillator. J. Appl. Phys. 2018, 124, 152132. 3.1. RQ1—Research Challenges in Edge Intelligence (EI)

3. Rafaie, M.; Hasan, M.H.; Alsaleem, F.M. Neuromorphic MEMS sensor network. Appl. Phys. Lett.

In the 1961 section, the 356 archers summarize the challenges faced by the Edge Intelligence (EI) paradigm that the

analyzed studies either mentioned or aimed to tackle. The discussion presented in this section aims to provide 4. Hasan, M.H.; Al-Ramini, A.; Abdel-Rahman, E.; Jafari, R.; Alsaleem, F. Colocalized Sensing and answers to RQ1: What are the main challenges and open issues in the distributed learning field? Intelligent Computing in Micro-Sensors. Sensors 2020, 20, 6346.

As Wentigned; 44dier, he Wanigg OIL Zendoj Qes Chere, edg Offene Wetwoed genais antellige of seing antomiteits, but it reidge corresoluting location and an still to rite beginning to IEEE Newler 2019, 35 till 5 being investigated. The surveyed studies tackle several challenges, which can be broadly grouped into six categories, displayed in Table 2 and described in what follows.

displayed in **Table 2** and described in what follows. 6. Li, E.; Zhou, Z.; Chen, X. Edge intelligence: On-demand deep learning model co-inference with

device-edge synergy. In Proceedings of the 2018 Workshop on Mobile Edge Communications, **Table 2.** Challenges in distributed machine learning in edge computing. Budapest, Hungary, 20 August 2018; pp. 31–36.

	Challenges	nd
CH1	Running ML/DL on devices with limited resources	ons for
CH2	Ensuring energy efficiency without compromising the accuracy	nent
CH3	Communication efficiency	94,

1		Challenges	
1	CH4	Ensuring data privacy and security	nance tional
	CH5	Handling failure in edge devices	. 2407–
1	CH6	Heterogeneity and low quality of data	ervice h 201_9 ittle
pr	016553412010	apally mainly when compared to the powerful data centers at the cloud. On the other hand,	many ML
13 ec re fo or 1 4	oplications Al-Rakh dge device Intellige quire stori Internat (DASC/ riginal train naMenge,is	require high computational power that outweighs the possibilities of resource-constrained ami, M.; Alsahli, M.; Hassan, M.M.; Alamri, A.; Guerrieri, A.; Fortino, G. Cost efficie s. Limited resources also include memory and storage capacities. NN and ML algorithms nce framework using docker containers. In Proceedings of the 2018 IEEE 16th ng of and access to a handful of parameters that describe the model architecture and weig ional Conference on Dependable, Autonomic and Secure Computing Congress classification model. With limited storage, it may not be possible to have continued acce PICom/DataCom/CyberSciTech), Athens, Greece, 12–15 August 2018; pp. 800–80 hing data, or the data may have been removed altogether to free up space. Therefore, a returning memory and storage, it may not be possible to have continued acce hing data, or the data may have been removed altogether to free up space. Therefore, a hing data, or the data may have been removed altogether to free up space. Therefore, a hing data, or the data may have been removed altogether to free up space. Therefore, a hing data, or the data may have been removed altogether to free up space. Therefore, a hing data, or the data may have been removed altogether to free up space. Therefore, a hing data is a comprehensive survey. IEEE Commun. Surv. Tutorials 2020, 22, 860	d IoT and ent edge generally ght values ess to the 07. significant dtiel@lernal
m	121990199 41990	Wiles ining. A comprehensive survey. The Commun. Surv. Tutonais 2020, 22, 805	9-904.
15 Cl In	. Zhou, Z H2 consist artificial general, t	.; Chen, X.; Li, E.; Zeng, L.; Luo, K.; Zhang, J. Edge intelligence: Paving the last m s of ensuring the energy efficiency of edge devices without compromising the accuracy of th intelligence with edge computing. Proc. IEEE 2019, 107, 1738–1762. he higher the complexity of the required processing, the more energy is consumed. Edge de	nile of le system. evices can
16	e Vatteraq	kæered; Waltusa, dasek, atzevederky oppren huptoon Jot Vagionalians Thiustaberinnaviazed. & Aæ	ehrvense igis
ef	fielisteiput	edever, thises heard in a doce with care to sur no (cosumpression for the data gene	erated and
1 th	e decision	s/inferences made. So there is an important trade off to be managed arning iot data patter	rns at
Cl pc Cl	the edg H3 concer Forum (oor connec (RTSI), onti [17] cla	e with hypothesis transfer learning. In Proceedings of the 2016 IEEE 2nd Internations communication issues, where edge intelligence models must consider that the devices room Research and Technologies for Society and Industry Leveraging a better tomorrow trivity. In such cases, the model update time in training tasks may be delayed. Valerio, Pass Bologna, Italy, 7–9 September 2016; pp. 1–6.	onal night face row arella and Illenges in
1&	onArAUniZat	idinin Xuuddevedlainenarfic; huuanatixnsHeikkwaanwidtin Wastariluund oTu Falgaaaned fogi soom	puting
	enabled	AI for IoT-an overview. In Proceedings of the 2019 IEEE International Conference	e on
C	H Aitsi fielia t	ebh te kligten na iv a Cir cuvids acudi (SyStevenal (AphOl AB)), n Si si nachge, i Mæivigæm, c & & a 2014/bænshi 20 (1915	appu (511.as
he	eatticare.	Thus, distributed ML algorithms must be able to preserve user privacy and information secu	urity when
da 19 pc	ata are trai . Hossair ossible ma edge co	nsferred throughout the devices. Distributed Edge-Intelligence (EI) has multiple points of vulne , M.S.; Muhammad, G.; Amin, S.U. Improving consumer satisfaction in smart citie licious attacks or leakage of confidential or important data in the ML workflow. mputing and caching: A case study of date fruits classification. Future Gener. Con	erability to s using nput.
CI	Syst. 20	118, 88, 333–341. challenge posed by failures in edge devices. Since devices might fail at some point, the c	distributed
2a	gShlanma	us constitution of the state and the state and the state of the state	elayineervolaura
ris	sePatadleea	ingsoffbien2021.8/Indeportionsale Departelyeince per Asekamaesine Compatiniah Command	eateionds on
th	e Chightroph	aling Netwoinginga (&CAOCEVA), the exter Noved and in diaply 21-200 citoty lige 200 8; oppar585-	15990 e the
20 21	ollected da Wan S	ta are sparse and unlabelled $\begin{bmatrix} 10 \\ . \end{bmatrix}$. Distributed edge intelligence can handle data from different to the edge for t	nt sources
	healthc	are service robots. Comput. Commun. 2020, 149, 99–106	

22. defendent by markets and este conjectory bioisectific Reprice align Sense reak or the network of the method of the sense state used and the sense of the sens

23. Rausch, T.; Dustdar, S. Edge intelligence: The convergence of humans, things, and ai. In Table 3 presents references to each of the described challenges, as well as studies that propose approaches to Proceedings of the 2019 IEEE International Conference on Cloud Engineering (IC2E), Prague, tackle these challenges. This table aims to only show an overview on the number of papers by each challenge. The Czech Republic, 24–27 June 2019; pp. 86–96.
researchers can observe that challenge CH1 is the one with more papers present in literature. All of the cited works
24re Fraseirates; Cibed/ituda nor, eF. Artificial Intelligence on Edge Computing: A Healthcare Scenario in Ambient Assisted Living. In Proceedings of the Artificial Intelligence for Ambient Assisted Living (AI*AAL.it 2019), Rended Ray, 2019; Norde Artificial Intelligence.

2		References	Works That Tackle the Challenges	n t. Pract.
2	CH1	[<u>10][15][18][19][20][21][22][23][24]</u> [<u>25][26][27]</u>	[14][24][27][28][29][30][31][32][33][34][35][36][37][38][39][40][41][42][43][44][45][46][47][48][49][50] [51][52][53][54][55][56][57][58][59][60][61][62]	hings
2	CH2	[10][15][19][22][24][26][37]	[8][20][24][27][29][32][33][34][45][48][52][55][57][61][62][63]	neous
	CH3	[10][20][24][25][42][62]	[16][17][19][29][30][31][32][35][39][42][52][64][65]	e on
2	CH4	[10][20][23][32][40][66]	[8][9][10][20][40][47][67][68][69]	nputing), 1–20.
2	CH5	[10][23]	_	M. Igs of Ington.
	CH6	[10][20][40][70][71]	[<u>10][34][66]</u>	
3	·· -·, -	-,	onon, /. Eugo / on domand docoloraling doop notice notice	(

inference via edge computing. IEEE Trans. Wirel. Commun. 2019, 19, 447-457.

31.2.i, H. Diang M. Learning IoT in edge: Deep learning for the Internet of Things with edge computing. IEEE Netw. 2018, 32, 96–101.

Here, the researchers focus on three main aspects namely: (i) the system architecture. (ii) how the ML tasks are 32. Hassan, M.A.; Xiao, M.; Wei, Q.; Chen, S. Help your mobile applications with fog computing. In distributed among the devices, and (iii) the underlying adopted techniques. The researchers classify the several Proceedings of the 2015 12th Annual IEEE International Conference on Sensing, Communication, approaches used in distributed learning based on these three aspects. The researchers identified nine groups of and Networking-Workshops (SECON Workshops), Seattle, WA, USA, 22–25 June 2015; pp. 1–6. techniques and strategies, described in what follows: Federated learning; Model partitioning; Right-sizing; Edge 331eLibbcessing; Scheduling; Cilchene Gaining; Kedge only; Multer Compression; and Sther Rechniques.

3.3nfRQB+ctFrankeworksforeEdgenintelligende, 249-261.

34. Liu, Y.; Yang, C.; Jiang, L.; Xie, S.; Zhang, Y. Intelligent edge computing for IoT-based energy This section describes the studies that provided answers to the RQ3 of this survey. Table 4 lists the main management in smart cities. IEEE Netw. 2019, 33, 111–117. frameworks currently used in distributed ML applications. The table also correlates each framework with the 35orNisphinding Mometanif Recipitentsselentionaiforefeaterated by arning with heterogeneous resources in mobile edge. In Proceedings of the ICC 2019-2019 IEEE International Conference on Communications (ICC), Shanghai, Chinkle 20 E244May 2019; pp. 1–7.

3	- ·	Groups of) national
(r	Framework	Techniques or	Comments	17–24.
J		Strategies		.rtIoT),
(T)	Neurosurgeon ^[72]	Model Partitioning	Lightweight scheduler to automatically partition DNN computation between edge devices and cloud at the granularity of NN layers	n of 2019,
З	JointDNN ^[73]	Model Partitioning	JointDNN provides an energy- and performance-efficient method of querying some layers on the mobile device and some layers on the cloud server.	stem 1962–
4	H. Li et al. ^[31]	Model Partitioning	They divide the NN layers and deploy the part with the lower ones (closer to the input) into edge servers and the part with higher layers (closer to the output) into the cloud for offloading processing. They also propose an offline and an online algorithm that schedules tasks in Edge servers.	ion for IN- 9, 6,
4	Musical chair ^[74]	Model Partitioning	Musical Chair aims at alleviating the compute cost and overcoming the resource barrier by distributing their computation: data parallelism and model parallelism.	t
4	AAIoT ^[75]	Model Partitioning	Accurate segmenting NNs under multi-layer IoT architectures	work for obal). 1–6.
4	MobileNet ^[42]	Model Compression Model	Presented by Google Inc., the two hyperparameters introduced allow the model builder to choose the right sized model for the specific application.	bia, J. et of
-1		Selector		. i i Ont.

 Guo, R.; Xiang, Y.; Mao, Z.; Yi, Z.; Zhao, X.; Shi, D. Artificial Intelligence Enabled Online Nonintrusive Load Monitoring Embedded in Smart Plugs. In Proceedings of the International Symposium on Signal Processing and Intelligent Recognition Systems, Trivandrum, India, 18–21 December 2019; Springer: Singapore, 2019; pp. 23–36.

4		Groups of		in the
Δ	Framework	Techniques or	Comments	8–64.
		Strategies		EE 29th
4	Squeezenet	Model Compression	It is a reduced DNN that achieves AlexNet-level accuracy with 50 times fewer parameters	vork
5	Tiny-YOLO	Model Compression	Tiny Yolo is a very lite NN and is hence suitable for running on edge devices. It has an accuracy that is comparable to the standard AlexNet for small class numbers but is much faster.	f cnn. In າ
5	BranchyNet	Right sizing	Open source DNN training framework that supports the early-exit mechanism.	S.
5	TeamNet ^[76]	Model Compression Transfer Learning	TeamNet trains shallower models using the similar but downsized architecture of a given SOTA (state of the art) deep model. The master node compares its uncertainty with the worker's and selects the one with the least uncertainty as to the final result.	y in Jn. or fog-
5	OpenEl ^[42]	Model Compression Data Quantization Model Selector	The algorithms are optimized by compressing the size of the model, quantizing the weight. The model selector will choose the most suitable model based on the developer's requirement (the default is accuracy) and the current computing resource.	–50. vork hou,
5	TensorFlow Lite ^[77]	Data Quantization	TensorFlow's lightweight solution, which is designed for mobile and edge devices. It leverages many optimization techniques, including quantized kernels, to reduce the latency.	nobile- is and

- 57. Gonzalez-Guerrero, P.; Tracy II, T.; Guo, X.; Stan, M.R. Towards low-power random forest using asynchronous computing with streams. In Proceedings of the 2019 Tenth International Green and Sustainable Computing Conference (IGSC), Alexandria, VA, USA, 21–24 October 2019; pp. 1–5.
- 58. Sanchez, J.; Soltani, N.; Chamarthi, R.; Sawant, A.; Tabkhi, H. A novel 1d-convolution accelerator for low-power real-time cnn processing on the edge. In Proceedings of the 2018 IEEE High

		Groups of		nber
5	Framework	Techniques or Strategies	Comments	⁻ eature ta
6	QNNPACK (Quantized Neural Networks PACKage) ^[78]	Data Quantization	Developed by Facebook, is a mobile-optimized library for high- performance NN inference. It provides an implementation of common NN operators on quantized 8-bit tensors.	۲iv
6	ProtoNN ^[79]	Model Compression	Inspired by k-Nearest Neighbor (KNN) and could be deployed on the edges with limited storage and computational power.	ecision tell.
6	EMI-RNN ^[80]	Right Sizing	It requires 72 times less computation than standard Long Short term Memory Networks (LSTM) and improves its accuracy by 1%.	d ervice- ·3385.
6	CoreML ^[81]	Model Compression Data Quantization	Published by Apple, it is a deep learning package optimized for on-device performance to minimize memory footprint and power consumption. Users are allowed to integrate the trained machine learning model into Apple products, such as Siri, Camera, and QuickType.	system nce on L1. system
6	DroNet ^[33]	Model Compression Data Quantization	The DroNet topology was inspired by residual networks and was reduced in size to minimize the bare image processing time (inference). The numerical representation of weights and activations reduces from the native one, 32-bit floating-point (Float32), down to a 16-bit fixed point one (Fixed16).	e. In ington, 17, 4,
6 6	Stratum ^[82]	Model Selector Dynamic Scheduling	Stratum can select the best model by evaluating a series of user- built models. A resource monitoring framework within Stratum keeps track of resource utilization and is responsible for triggering actions to elastically scale resources and migrate tasks, as needed, to meet the ML workflow's Quality of Services (QoS). ML modules can be placed on the edge of the Cloud layer, depending on user requirements and capacity analysis.	og-to- Is. IEEE Ibject), 7,
	3681-3692.			

70. Zhou, Z.; Liao, H.; Gu, B.; Huq, K.M.S.; Mumtaz, S.; Rodriguez, J. Robust mobile crowd sensing: When deep learning meets edge computing. IEEE Netw. 2018, 32, 54–60.

7		Groups of		elligent
7	Framework	Techniques or	Comments	
		Strategies		rchit.
7	Efficient distributed deep learning (EDDL)	Model Compression	A systematic and structured scheme based on balanced incomplete block design (BIBD) used in situations where the dataflows in DNNs are sparse. Vertical and horizontal model	engine -576.
7		Model Partitioning	partition and grouped convolution techniques are used to reduce computation and memory. To speed up the inference, BranchyNet	
7		Right-Sizing	is utilized.	ms.
7	In-Edge AI ^[5]	Federated Learning	Utilizes the collaboration among devices and edge nodes to exchange the learning parameters for better training and inference of the models.	n /stems
7	Edgence ^[83]	Blockchain	Edgence (EDGe + intelligENCE) is proposed to serve as a blockchain-enabled edge-computing platform to intelligently manage massive decentralized applications in IoT use cases.	online: Inpack
7	FederatedAveraging (FedAvg) ^{[<u>84]</u>}	Federated Learning	Combines local stochastic gradient descent (SGD) on each client with a server that performs model averaging.	Jdupa, es. In
8	SSGD ^[85]	Federated Learning	System that enables multiple parties to jointly learn an accurate neural network model for a given objective without sharing their input datasets.	1 nt
8	BlockFL ^[86]	Blockchain Federated Learning	Mobile devices' local model updates are exchanged and verified by leveraging blockchain.	
8				erless
	Edgent ^{[<u>6]</u>}	Model Partitioning	Adaptively partitions DNN computation between the device and edge, in order to leverage hybrid computation resources in	n 9),

83. Xu, J.; Wang, S.; Zhou, A.; Yang, F. Edgence: A blockchain-enabled edge-computing platform for intelligent IoT-based dApps. China Commun. 2020, 17, 78–87.

8		Groups of		nt
	Framework	Techniques or	Comments	PMLR:
8		Strategies		CM
8		Right-Sizing	proximity for real-time DNN inference. DNN right-sizing accelerates DNN inference through the early exit at a proper intermediate DNN layer to further reduce the computation latency.	10
8	PipeDream ^[87]	Model Partitioning	PipeDream keeps all available GPUs productive by systematically partitioning DNN layers among them to balance work and minimize communication.	ns, P.
8	GoSGD ^[88]	Gossip Averaging	Method to share information between different threads based on gossip algorithms and showing good consensus convergence properties.	using Alere are They are: H., Kim;
	Gossiping SGD ^[89]	Gossip Averaging	Asynchronous method that replaces the all-reduce collective operation of synchronous training with a gossip aggregation algorithm.	Federated eep less than sum on
5	GossipGraD ^[90]	Gossip Averaging	Asynchronous communication of gradients for further reducing the communication cost.	1 on
Ĝ	INCEPTIONN ^[91]	Data Quantization	Lossy-compression algorithm for floating-point gradients. The framework reduces the communication time by 70.9 80.7% and offers 2.2 3.1× speedup over the conventional training system while achieving the same level of accuracy.	mobile າy, 10–
G	Minerva ^[92]	Data Quantization Model	Quantization analysis minimizes bit widths without exceeding a strict prediction error bound. Compared to a 16-bit fixed-point baseline, Minerva reduces power consumption by 1.5×. Minerva identifies operands that are close to zero and removes them from	ıre deep mational 2018;
9! 9!	Cloud Computing, C 5. Sun, W.; Liu, J.; Yue meets industrial IoT.	arisbad, CA, U 0% , Y. Al-enhanc IEEE Netw. 2	USA, 11–13 October 2018; pp. 401–411. 5% 10% 15% 20% 25% ced offloading in edge computing: When machine learning 2019, 33, 68–74.	work m on 30% ng

g		Groups of		ıstry
g	Framework	Techniques or	Comments	nd testing idustrial
		Strategies		ven to be
			affected. Selective pruning further reduces power consumption by	terms of
g			$2.0 \times$ on top of bit width quantization.	rtitioning).
	AdaDeep ^[93]	Model	Automatically selects a combination of compression techniques	tea ownhole
	·	Compression	for a given DNN that will lead to an optimal balance between user-	28–236.
G			specified performance goals and resource constraints. AdaDeep enables up to 9.8× latency reduction, 4.3× energy efficiency	In using data (ICCE- g security
			negligible accuracy loss.	nting their
10				nce-
	JALAD ^[94]	Data Quantization	Data compression by jointly considering compression rate and model accuracy. A latency-aware deep decoupling strategy to	13, 10,
10		Model	minimize the overall execution latency is employed. Decouples a	felððf El
10		Partitioning	inside the conventional cloud.	ossible to () 바ealth,
		0 <u>201</u> .	Ale for the second state of the second state of the state of the second state of the s	research.
10 ₂	3 Liu L.: Zhang X.: Z summarizes the statistics safety service on ve	hang, O.: Wei of the six domai hicles. In Proc	nert, A.: Wang, Y.: Shi, W. AutoVAPS: An IoT-enabled p ns of the publishing by field. ceedings of the Fourth Workshop on International Scien	. Figure 3 ublic ce of
	Smart City Operatio	ns and Platfor	ms En gineering, Montr eal, QC, Canada, 15 April 2019;	рр. 41—
	47.		Domain Area	
10	 Yang, H.; Wen, J.; V multipedestrian trac 	Vu, X.; He, L.; king method w	Mumtaz, S. An efficient edge artificial intelligence	1178_
	4188.			
10	5. Lee, C.; Park, S.; Ya Industry sensing based on e	ang, T.; Lee, S ^{eillance} dge computing	.H. Smart parking with fine-grained localization and use Security g. In Proceedings of the 2019 IEEE 90th Vehicular Tech	r _{energy} us
10	Conference (VTC20 Industry 4:0 Ce (Analytics and Smart Industry Envir Envir	19-Fall), HONO Data Analysis Cyt onment	Diulu, HI, USA, 22-25 September 2019; pp. 1-5. Public Transport Emotion recognition	lanagement Smart Grid Smart City
τU	Industry Control Min S., Mage R IoT-Network fusion for a cognitive	ewireless fran	Blockchain Nework. IEEE Wirel. Commun. 2019 Honder Healthcare	mart Home

107. Cao, Y.; Hou, P.; Brown, D.; Wang, J.; Chen, S. Distributed analytics and edge intelligence: Pervasive health monitoring at the interval fegatoric fega

- 108. Al-Rakhami, M.; Gumaei, A.; Alsahli, M.; Hassan, M.M.; Alamri, A.; Guerrieri, A.; Fortino, G. A lightweight and cost effective edge intelligence architecture based on containerization technology. World Wide Web 2020-23-1341-1360
- 109. Rincon, Surveillance, V.; Carrascosa, C. Towards the Edge Intelligence: Robot Assistant for the Detection and Classification of Human Emotions. In Proceedings of the International Conference on Practical Applications of Agents and Multi-Agent Systems, L'Aquila, Italy, 13–15 July 2020; Springer: Beffin/Heidelberg, Germany, 2020; pp. 31–41.
- 110. Maitra, A.; Kuntagod, N. A novel mobile application to assist maternal health workers in rural India. In Proceedings of the 2013 5th International Workshop on Software Engineering in Health Care (SEHC), San Francisco, CA, USA, 20–21 May 2013, pp. 75–78.
- 111. Gupta, A.; M바위바jee, N. A cloudlet platform with virtual sensors for smart edge computing. IEEE Internet Things J. 2019, 0, 0435—0402.
- 112. Energy Managementart home system architecture facilitated with distributed and embedded flexible edge analytics in demand-side management. Int. Trans. Electr. Energy Syst. 2019, 29, e12014. Retrieved from https://www.@ncyclopedia.pub/entry/history/show/52931 10 15

Figure 3. Publications by domain application.

Table 5. Application domains and corresponding works.

Domains	Works That Approach the Theme
Industry (8)	[8][27][47][49][64][95][96][97]
Surveillance (5)	[38][65][98][99][100]
Security (4)	[<u>9][25][67][101]</u>
Intelligent Transport Systems (ITS) (13)	[22][36][39][41][47][48][52][54][71][102][103][104][105],
Health (14)	[13][21][24][43][44][47][63][66][106][107][108][109][110][111]
Energy Management (4)	[34][46][47][112]