# Landing AI on Networks: An equipment vendor viewpoint on Autonomous Driving Networks

Dario Rossi[1], Liang Zhang[2]

[1]*Huawei Technologies, co. Ltd, Paris Research Center* — dario.rossi@huawei.com
[2]*Huawei Technologies, co. Ltd, Nanjing Research and Development Center* — zhangliang1@huawei.com

*Abstract*—The tremendous achievements of Artificial Intelligence (AI) in computer vision, natural language processing, games and robotics, has extended the reach of the AI hype to other fields: in telecommunication networks, the long term vision is to let AI fully manage, and autonomously drive, all aspects of network operation. In this industry vision paper, we discuss challenges and opportunities of Autonomous Driving Network (ADN) driven by AI technologies. To understand how AI can be successfully landed in current and future networks, we start by outlining challenges that are specific to the networking domain, putting them in perspective with advances that AI has achieved in other fields. We then present a system view, clarifying how AI can be fitted in the network architecture. We finally discuss current achievements as well as future promises of AI in networks, mentioning a roadmap to avoid bumps in the road that leads to true large-scale deployment of AI technologies in networks.

*Index Terms*—Artificial intelligence; Machine Learning; Network Management; Network O&M; AI-Native;

## I. THE NEW GOLD

**T**HE last decade has witnessed significant advances in several fields where Artificial Intelligence (AI) has been applied to – from image recognition, to natural language processing and gaming to name a few examples. Such achievements are due to the fortunate confluence of several necessary ingredients: namely, (i) exceptional theoretical advances in the last 50 years, coupled to the availability of (ii) massive volumes of data, and equivalently (iii) massive computing capabilities. These achievements have gained significant press attention, fueling the hype on AI techniques, and their expected benefits. As a result, every technology sector joined this new "Gold rush", including the networking field, where AI is envisioned on the long-run to fully manage, and autonomously drive, all aspects of network operation [1].

At the same time, a significant fraction of AI projects are difficult to transfer[1] beyond the initial proof-of-concept. The difficulty in successfully landing AI is overtly recognized lately [3], and is technically rooted in either the lack of some of the above ingredients, difficulty in integration, or other non-technical aspects. As "not all gold that glitters," it is necessary to understand what type of problems AI can solve, and how AI solutions can be fit in the overall system: this is necessary, in order for AI to really make the difference in a specific technology field, such as the networking domain considered in this paper, before the next AI winter.

A decade ago, Marc Andreessen was rightly anticipating that "software is eating the world": in the last decade, evolution toward software enabled networking world to escape ossification [4], and AI-software shows the very same appetite for the next decade. Owing to growing success of all-IP (2000-2010) and cloud-native (2010-2020) networking technologies, IP-enabled communications are now spanning a very large (and still growing) set of vertical sectors and markets. To manage such a plethora of heterogeneous services evolving at a fast pace, the network operation and management (O&M) community has started turning its attention to AI, for relieving and assisting human operation for diverse tasks (e.g., ranging from configuration, to dynamic resource management, troubleshooting and quality assessment). In a field where a significant fraction of the operations are still involving human intervention, and where such interventions are also responsible for a significant fraction of the errors, AI seems an appealing means to immediately automate part of these manual tasks (e.g., from configuration, to fine-grained resource management at very fast timescale, troubleshooting guidance and quality forecast) and later reach fully automated and error-free (or at least self-healing) operations.

Clearly, the evolution of the network O&M to a fully autonomous driving network (ADN), cannot be done overnight due to technical challenges, adoption barriers and legal aspects (e.g., liability). Pragmatically, our vision is thus for AI to replace human hands in the *fast loop*, but not fully supplant humans which is essential to keep in the *slow loop*. Taking the viewpoint of an equipment vendor, this paper illustrates the current status of network AI, enriching the narrative with examples of research that successfully landed in network technologies, further laying out the steps necessary to the ADN for reaching higher level of autonomicity and intelligence. We instead disregard complementary technological aspects (well covered by [5], [6]) and methodological aspects (for which we refer the reader to [6]–[10]).

The rest of this paper is organized as follows. Sec. II briefly introduces AI, and Sec. III overviews AI challenges on the network domain, putting them in perspective against other fields where AI has been successful. Sec. IV then introduces the key aspect of the Autonomous Driving Network (ADN), examine its architectural, hardware and software needs with a focus on AI-related aspects. Finally, Sec. V overviews illustrative examples of how AI can be successfully landed in current networks, while Sec. VI discusses open and future challenges on the path towards the ADN.

---

[1]As Gartner put it [2], in 2020 "80% of AI projects will remain alchemy, run by wizards whose talents will not scale in the organization"

## II. WHAT IS AI, AND WHY IT MATTERS ?

As AI has recently become an abused term, and given that the very same definition of AI is constantly redefined[2], we briefly introduce it here what we consider to be the set of AI techniques that will constitute ADN's basic building blocks.

### A. Brief history of AI & ML

As visually depicted in the top part of Fig.1, the history of AI & ML can be traced traced back to Alang Turing's work in the early 1950s with the basis of Neural Networks (NN), as well as the terms Artificial Intelligence (AI) and Machine Learning (ML) introduced during that decade[3]. Clearly, AI/ML techniques have evolved significantly since then: as for many scientific fields, their evolution has been shaped by "hype" cycles, where peaks of attention and spending sprees (known as "AI springs") were followed by periods of disengagement and funding cuts (known as "AI winters", in the late 70s and late 80s). While the network community pioneered the use of AI as early as in the 70s [13], many of the now popular branches of AI, such as Deep Learning (DL) and Data Mining were introduced in the last AI spring[4].

### B. Crude taxonomy of AI & ML

As AI is again reaching very high hype levels, it is a valid question to ask which AI techniques, and to what extent, can be successfully landed in the network field before the next "ice age" of AI. In the context of this paper, we are not interested in establishing a rigorous taxonomy of AI & ML techniques. Still, given that these terms have been abused up to being the target of popular jokes [15], it is useful for clarity (and simplicity) to refer to a (crude) taxonomy reported in the bottom of Fig. 1. For the scope of this paper, we restrain the focus to techniques that belong to the data-driven branch of AI that is also known as ML. By abuse of language, we will use the two acronyms interchangeably in the following.

From a very high level, the purpose of ML can be either to (i) make effective use of *existing* knowledge, or (ii) gather structured understanding on *unknown* phenomena, as well as (iii) *learn* to achieve a goal. Roughly, these purposes directly map to three branches of ML that are respectively known as (i) supervised, (ii) unsupervised and (iii) reinforcement learning – though blending among categories is possible (e.g., semi-supervised or self-supervised learning). While a more detailed and comprehensive taxonomy is out of scope, examples of relevant techniques for each branch are provided in the bottom of Fig. 1. Oversimplifying, many (though not all) of the ML techniques are efficient ways to solve complex data-driven
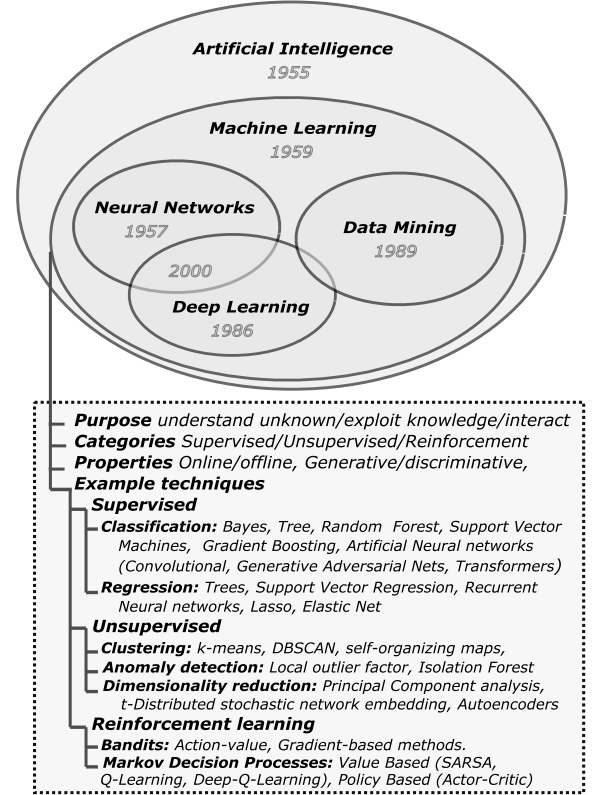
Fig. 1. Brief history and crude taxonomy of AI.

optimization problems, offering solutions that are well suited for the data at hand (i.e., fit well), but are also amenable to generalization (i.e., avoid overfitting). This constitutes part of the reasons of ML success in several fields, and makes it practically appealing for the network domain too.

### C. Example of AI success across all domains

The most recently hyped examples of success in the current "AI spring" pertain to areas such as computer vision, game-playing and natural language processing.

Image recognition attracted significant attention not only as being among the first key success of Convolutional Neural Networks (CNN) [14], but also e.g., in reason of the powerful Generative Adversarial Networks (GAN) [16] underneath the (in)famous DeepFakes technologies [17]. In the game-playing context, Deep Reinforcement Learning (DRL) has been instrumental in achieving super-human playing abilities, with e.g., Google's AlphaGo [18] beating the go board game world champion Lee Sedol, or OpenAI Five [19] winning the online computer-game DOTA2 tournament. In natural language processing, self-supervised neural embedding (e.g., *word2vec* [20]) and few-shot transformer technologies [21] such as OpenAI GPT3 [22] gained significant traction lately.

As we illustrate in Sec.V, communications and computer networks are one among the numerous other domain of applications (biology, medical field, robotics) that is currently exploring the use of AI techniques in many aspects of its operation.

## III. Landing AI on Networks

To replicate AI success in other fields, it is necessary to first understand their root cause: we thus start by dissecting AI successes, to draw informed conclusions from a networking perspective[5] for the ADN.

Without willing to undermine the significance of AI achievements, we remark that recent advances have equally benefited from: (i) ground-breaking advances in ML theoretic research, (ii) the availability of large corpora of (labeled) data to feed ML models with, and (iii) the availability of software platforms that efficiently exploit hardware acceleration. Two observations are worth sharing: on the one hand, whereas the latest wave of AI success essentially involves neural networks, this is clearly not the first time that AI accomplishes similar prowess; on the other hand, this AI spring may be the first time at which time is ripe for advances in all (i), (ii) and (iii) aspects. Aside from the above technical aspects, success also depends on (iv) business considerations. These four aspects are crucial for landing AI in networks, as we next discuss.

### A. Theoretical advances

**Root of AI success**. The latest AI spring yielded to numerous theoretic advances. The most visible ones, in reason of their recent hype, involve theoretic advances of Deep Neural Networks [23], for which Y. Bengio, Y. LeCun and G. Hinton were recently distinguished with the 2018 ACM Turing award. Other equally important advances include such as the ensembling of several "weak" classifiers – that became widely known after the popular Netflix-prize [24]. Ensembles are now state of the art for many supervised [25], [26] and unsupervised [27] tasks. Similarly, advances in causal reasoning capabilities [28] for which J. Pearl has been credited the 2011 ACM Turing award, are perhaps less known to the broad public, but equally relevant. Finally, more prospective advances (e.g., spiking neural networks [29]) are under way, but are in still early stage of development so they can be expected to have a deeper impact on a longer time horizon.

**Landing AI in Networks**. Fortunately, the greatest inventions[6] of our time has democratized access to knowledge: as such, all AI advances can be readily used in the network field. Clearly, whereas many of these recently hyped successes involve DL technologies, it appears evident that also any of the lesser-hyper AI technologies reported in Fig.1 is still worth considering from the viewpoint of an AI-fueled autonomous driving network. Purposely, we later (see Sec.V) report on successful network O&M application using AI techniques spanning all AI branches.

---

[5]This paper adopts an "AI for networks" angle, i.e., how AI can improve network Operation & Management (O&M). We point out that, in reason of AI success, an equally important viewpoint that research has considered is "networks for AI", i.e., how networking techniques can make AI workflow more efficient, where AI becomes thus a networked application – which we instead disregard due to space constraints.

[6]The foundation of the TCP/IP architectural principles of the Internet is traced back to mid 70s, for which V. Cerf and B. Kahn were credited with the 2004 ACM Turing award. Similarly, the foundation of Internet's most successful application, i.e., the World Wide Web, is traced back in the late 80s, for which T. Berners-Lee received the 2017 ACM Turing award.

TABLE I
Generic CPU, GPU vs Domain-specific hardware accelerators

| Vendor | Product | Target | Processing [TOPS] | [TFLOPS] | Power [W] |
|--------|---------|--------|--------|---------|-------|
| ARM | Cortex A72 [40] | Edge CPU | n.a. | 0.03 | 0.75 |
| Google | Coral.AI [41] | Edge DSA | 4 | n.a. | 2 |
| Huawei | Ascend310 [42] | Edge DSA | 22 | 11 | 8 |
| Intel | Xeon 8280 [?] | Cloud CPU | n.a. | $\approx 2$ | 205 |
| NVIDIA | P100 [43] | Cloud GPU | n.a. | 5-21 | 250 |
| Google | TPUv3 [44] | Cloud DSA | n.a. | 420 | $\approx 300$ |
| Huawei | Ascend910 [45] | Cloud DSA | 640 | 320 | 310 |

### B. Hardware and software

**Root of AI success**. Ultimately, theoretical advances are distilled in algorithms, for which hardware and software engines running them are equally important. From a *hardware* viewpoint, DNN have benefited from the emergence and commoditization of hardware acceleration such as Graphic Processing Units (GPUs), that are widely acknowledged to have substantially contributed to the recent success of AI [30]. A further evolution is the emergence of domain-specific architectures (DSA), that are well discussed by J. L. Hennesey and D. Patterson, in their 2017 ACM Turing award inaugural lecture [31]: Tab. I reports examples of DSAs for NN acceleration. Worth mentioning are advances on binary [32] and xor [33] neural networks, as well as neuromorphic processors [34] for spiking networks acceleration, of interest on a longer time horizon.

The availability of *software* stacks able to efficiently leverage the above hardware in a seamless manner, jointly providing a unified and complete environment is key to lower bootstrap cost of AI in any field, as well as to facilitate transfer. Popular stacks include Google TensorFlow [35], Huawei MindSpore [36], Telsa Pytorch [37] for neural network, and Scikit-learn [38], for general workflows. These stacks offer the ability to rapidly prototype in high-level language, with Python commonplace now in the scientific workflow, and have optimized backends seamlessly supporting hardware acceleration. Software stacks have been at least as important as hardware accelerators (if not more according to [39]) to lower AI entering barrier, with respect to the more scattered situation of just less than a decade ago.

**Landing AI in Networks**. The set of hardware accelerators and software stacks just introduced for the general AI case, undoubtedly facilitates development and execution of AI workflow for network as well. At the same time, whereas Cloud-native workflows can expect to find GPUs or high-end TPUs of Tab.I, many of the network operations will need to be carried out without having access to cloud resources. Similarly, whereas the factor form and their power drain of some TPU chipsets (eg. Ascend310 or Coral.AI) are small, however this does increase Capex (for the new chip) and Opex (higher computational cost, albeit small, for the new AI feature). As such, whereas SmartNICs equipped with such powerful AI accelerators start to appear [46], it is difficult to predict yet how much they will be widespread in networking

equipments. Thus, a safer bet is to consider an AI workflow that can exploit AI acceleration if available, but is lean enough to run on standard CPUs (e.g., as for the case of ARM, using software acceleration libraries such as armNN [47], that can seamlessly support Cortex-A CPUs and Mali GPUs).

Additionally, to optimize communication between CPU and the GPU/TPU boards/chipsets, GPU and TPUs are designed for batch processing, with furthermore relatively large batch size of several thousands elements – which is in contrast with the needs of typical network workflow. Indeed, while networks stack employ batch processing for *packet-level* operation [48], the size of batches is one to two orders of magnitude smaller. Additionally, in typical networking use-cases, most AI processing should happen instantaneously, i.e., with batch size equal to 1. Considering for instance *per-flow* decisions, it appear inconvenient to batch AI processing across multiple flows: e.g., waiting for the arrival of several flows due to batching would require buffering traffic, couple decisions and delay application of AI-driven policies. Fortunately, the emergence of stacks such as TensorFlow-Lite (TFL) [49], or more recently TensorFlow-Lite micro (TFLM) [50], makes it easier to run proof-of-concept AI models on constrained devices and CPUs that are pervasive in network equipment.

Shortly, we argue that while hardware accelerations starts being available, due to the KISS principle, AI solutions that are excessively computationally costly are not going to be successfully landed in networks – an aspect worth further attention that we thus consider in the following discussion.

### C. Data (and environment)

**Root of AI success**. The third key of success has been the availability of large datasets (or controlled environments). For instance, MNIST [51], CIFAR [52] or ImageNet [53], that overall comprise tens of millions of images for tens of thousands of classes, have been instrumental to fuel advances in image recognition [14]. Similarly, recently hyped advances [22] in NLP relied on hundred billions of text tokens corpus such as Common Crawl [54]. Even in lesser mass-mediatized fields such as computational biology, it is fairly well recognized that astonishing advances [55] would not have been possible without 50 years of expert-driven labeling work on protein unfolding. In the reinforcement learning branch of AI, one or more agents interact with an environment to learn a successful strategy, by enforcing actions that alter environmental responses. In this playground, OpenAI Five [19] learned by playing over 10,000 years worth of games against itself and AlphaGo-Zero [18] was trained with 29 million games of self-play during 40 days using 4 TPUs. Clearly, the ability to super-scale the exploration of the action space in a faster-than-real-time yet realistic-enough environment has been crucial to achieve such results.

**Landing AI in Networks**. Unconstrained access to high quality data, which is key for accurate models with good generalization capabilities, is a known problem across all AI domains of applications [56], and networking is not an exception. Interaction with an environment in a closed learning loop further exacerbates the problem.

Considering the classic "4V" data properties for the sake of simplicity, in terms of *volume* and *velocity*, it is sufficient to recall that the data rate of a single ToR switch is higher than the high-volume physics collected by Stanfords' Large Synoptic Survey Telescope [57] or CERN's Large Hadron Collider (LHC) [58]. Additionally, if the widely popular Moore law postulated a exponential growth in the computing capacity, the lesser known but equally important Gilder law observes that "Bandwidth grows at least three times faster than the compute power", i.e., making the matter worse than in other fields. In terms of *variety*, network data is extremely scattered, multi-modal, heterogeneous, more than what typically happens in domains such as image or natural language processing. Therefore, a standard and homogeneous data representation would significantly facilitate AI application in networks. Furthermore, such heterogeneity make AI generalization capabilities of paramount importance.

Finally, in terms of *veracity*, it appears evident that whereas other fields have crowdsourced, amassed and shared large datasets, the networking field is lagging far behind. This is due legal/business constraints on data sharing on the one hand (see Sec.VI-A), and on the intrinsic difficulty of the labeling task on the other hand. For instance, for image or NLP use-cases industries are either crowdsourcing labeling to the whole population of Web users (e.g., where every day an estimated 500 years of manpower [59] is used to solve CAPTCHA puzzles to identify buses, cars and pedestrian in images, for training the future generation of self-driving vehicles), or directly recruiting human labor (e.g., speech recognition for digital personal assistants employs significant amount of human labor [60] at all steps, including to sanitize, verify and validate transcriptions). Conversely, labeling of, e.g., network anomalies, or identification of their root cause, or encrypted application traffic, is significantly more difficult to crowdsource as it requires domain expertise and large amount of time of skilled workers, which makes labeling cost higher.

On the one hand, we remark that a number of techniques and good practices can help with *data-related* (e.g., few-shot [61] or active [62] learning techniques both help reducing the number of labels) and *environment-related* (pre-train models for transfer learning and fine-tuning in the real-environment [63]) problems. On the other hand, we stress that AI solutions that are unable to be continuously updated and fail to seamlessly generalize are not going to be successfully landed in networks either – a second aspect worth to systematically consider in the remainder of this paper.

### D. Business (and beyond)

**Root of AI success**. Concurring in the ultimate success of a technology are also non-technical aspects, which are at least worth mentioning. For instance, for AI to be successful, it needs to solve an open (or otherwise unsolvable) problem, or solve it in a (significantly) more cost-effective way. By rephrasing Hockam razor, we may consider that if simpler solutions are available yielding good enough results at a fraction of the cost, then the simpler solution will be adopted (e.g., as it happened for the winner solution of the Netflix prize

[24], that was finally not deployed). Organizational aspects are complementary and equally relevant: these include, the relevance of the selected business case, the availability of AI profiles (data scientist, data engineers), and the culture of the enterprise (e.g., the amount of effort in collecting, sanitizing and treasuring the data). These aspects, which we will not discuss in this paper due to space limit, are surveyed e.g., in [64], [65].

**Landing AI in Networks**. The above business considerations are relevant for the networking domain too. While this paper mostly considers AI from a technical viewpoint, first discussing how AI fits in the network architecture (Sec.IV) and next introducing examples of AI usage in current networks (Sec.V), it will also devote space to discuss legal (Sec.VI-B) and human interaction (Sec.VI-C) aspects, that are particularly relevant for successful application of AI to networking.

## IV. ANATOMY OF AN AUTONOMOUS DRIVING NETWORK

While the specific domain of application of AI in networks are very diverse (e.g., from fiber/WiFi/5G/6G access, from campus to WAN and data-center networks to name a few) we can identify commonalities across all these different environment. Indeed, while the resources to manage and the goal of their management differ across the environments, we argue that the underlying family of AI algorithms to empower the current and future generation of networks and protocols are similar, can be fit in a single unified architecture. Our aim in this paper is not to present fully-fledged architectural details, which is the goals of standardization fora[7]. Rather, this sections aims at introducing the "DNA of the ADN", i.e., the broad *architectural principles* and the key hardware and software *building blocks* of the autonomous driving network.

### A. Three-tiered Network AI Architecture

The ADN comprises several physical elements, arranged from a logical point of view in a three-tiered architecture: this stem from the fact that spatial properties of the network, and the time constraints for its management, require a distributed and hierarchical organization of the AI functions. Spatio-temporal aspects are reflected in the synoptic illustrated in Fig. 2: the lowest tier is represented by fast and local AI decisions, whereas online AI actions requiring a network-wide knowledge fit the middle tier, and finally offline AI tasks with global multi-network significance fit the top tier. Open APIs are necessary for northbound interfaces, to interoperate with third party cloud platforms and AI software stacks: the top two levels of the architecture are also instrumental into automatically converting "intent" into configuration actions. Open APIs are also needed for southbound interfaces, to accommodate third party devices: the bottom levels of the architecture interact by upstreaming measurement from devices, as well as downstreaming configuration actions, in a closed ADN control loop.

[7]Relevant fora are e.g., ETSI Experiential Network Intelligence [66], ETSI Zero-Touch [67] and IRTF Computing in the Network (COIN) [68].
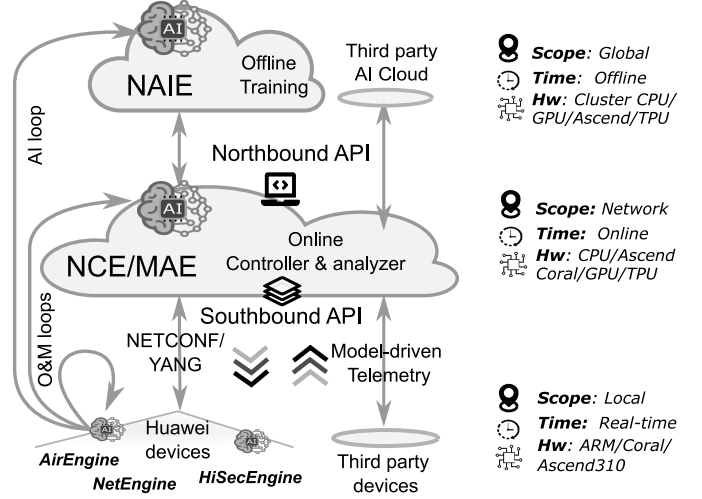


Fig. 2. *Network AI architecture*: three-tiered logical structure to fit different AI tasks from fast and local (device), to interactive and network-level (fog/cloud), to offline and global (cloud). The tiers interact in two kinds of closed loops for O&M and AI operations (see Sec. IV-B).

**Device AI**. Local decisions that have to be taken in near-real time are delegated close to or onboard of devices (e.g., Huawei's Net/Air/HiSec Engines [69]–[71]). For instance, augmented "visibility" tasks that require per-packet or per-flow decisions (scheduling, shaping, application identification, QoE estimation, anomaly detection) can be run at the edge. Devices (e.g., WLAN AP or datacenter switches) can take autonomous decisions (e.g., auto-configuration and tuning) based on their local knowledge. Such tasks might benefit from low-power hardware acceleration chipsets for AI tasks (recall Tab.I) or just be equipped with standard CPU and benefit from software acceleration, so that bespoke and highly-specific AI software stacks are expected to be the norm.

Further, non-technical considerations can favor in-device processing (e.g., business-sensitive or privacy-related for GDPR compliance) to keep data local to CPE vs processing data remotely. Other considerations, related to hardware availability or processing vs bandwidth cost (e.g., deploying many weak in-device accelerators vs fewer fog-accelerators or even fewer but more powerful cloud devices) can also affect the suitability of local vs remote processing for specific tasks.

**Online fog/cloud AI**. Decisions that require network-level knowledge, and that do not have strict sub-second latency requirements, such as for controller and analyzer tasks, can be offloaded to the fog/cloud (such as Huawei's iMaster NCE/MAE [72], [73] for the fixed/mobile network segments respectively). Controller actions can complement (or substitute) the one taken by devices, such as adding a slower centralized intelligence (e.g., taken by a single WLAN AP controller) on top of fast distributed decisions (e.g., taken by several WLAN APs). Analysis tasks can provide a broader knowledge than that locally accessible, by e.g., correlating anomalous events at network scale for troubleshooting.

From a hardware viewpoint, fog/cloud resources amortize Capex investments related to hardware acceleration, so that GPU and TPU should be expected to be more easily available whereas from a software perspective, code can leverage common and popular AI stacks. Controllers can access multi-vendor devices using cross-vendor southbound APIs (e.g., NetConf/Yang) to both upstream model-driven telemetry to feed AI decisions, as well as downstream automated configuration decisions to devices. From a business viewpoint, the fog/cloud AI can be offered as a service, which tradeoffs Capex investments for higher bandwidth usage.

**Offline cloud AI**. Knowledge that goes beyond the operation and management of a single network, or that span a large timescale, is better fit for offline storage and processing in the cloud (such as Huawei's iMaster NAIE [74]). This includes for instance model catalog management, model training services, transfer learning, federated learning, model health tracking, etc. These actions are instrumental for the good operation of network AI models, but are infrequent or anyway hardly need to be performed in a dynamic and interactive fashion (e.g., even if models are updated daily, data collection can be continuous, while training can be done nightly). As the offline AI cloud itself can become a bottleneck, research on the complementary "network for AI" viewpoint tackles the optimization of AI workloads from a network system perspective (see for instance [75], [76] and references therein).

From a hardware viewpoint, these AI tasks can leverage a large fleet of cloud resources, spanning several type of AI accelerators. From a software viewpoint, the offline cloud can offer managed services (e.g., Huawei's ModelArts [77]) but remains open and compatible with alternatives stacks (e.g., Amazon SageMaker [78]) and model marketplaces (e.g., Acumos [79]), further supporting open-source development through open APIs.

### B. Network AI and autonomy levels

With reference to autonomous driving vehicles, the industry identifies several levels of increased autonomy, from driver assisted (L1) to partial (L2), conditional (L3), high (L4) and full (L5) automation. A similar categorization has then be extended to the context of network automation: the goal of the ADN is to transition toward increasingly autonomous loops, where frequent human intervention (e.g., for technical necessity) is gradually substituted with sporadic human supervision (e.g., for legal aspects).

In the path that leads to a fully autonomous network, we can identify and map the needed AI and ML techniques to reach a given level. In simple terms, we differentiate between techniques that need to be applied in *open-loop* at that can greatly assist in augmenting the knowledge about the network operation (L2-3), as well as techniques applied in *closed-loop* to increasingly actuate or learn from the network (L4-5). We now briefly describe these simplified categories from the viewpoints of O&M and AI loops illustrated in Fig. 2.

**O&M loop**. An important aspect is to consider if the

AI-enabled O&M building block is working on open-vs closed-loop mode from a networking perspective. Supervised/unsupervised AI techniques are fit for *open-loop O&M* tasks such as application identification, traffic and quality forecast, imputation for missing data, compressed/enhanced telemetry, fault and anomaly detection etc. It should be clear that such techniques mainly needs data to be fueled (and possibly labels for supervised learning) and can be either operated in open-loop (sufficient for L2-3) or closed-loop AI modes (necessary at L4-5). These building blocks are fit for application on devices (or at the edge, depending on the timeliness vs computation requirements), although some tasks require the cloud processing power (for training or big data analytics).

*Closed-loop O&M* can leverage AI techniques for several tasks, ranging from resource allocation, configuration adaptation, fault prevention and repair. O&M loop can be fully distributed, or have a termination point in the fog/cloud, where centralized decisions can complement distributed decision. As learning directly from the real deployment can be hazardous (performance during a cold-start learning phases will be bad), it is desirable to pre-train in an actionable controlled environment (e.g., simulator, emulator) before further refine learning (e.g., digital twins, real network). AI techniques for closed-loop O&M are intrinsically operated in closed-loop AI mode, although the O&M loop can be closed in different network architectural points (i.e., device/edge/fog/cloud) depending on the specifics constraints of the application use-case (e.g., latency, telemetry, timescale, processing power, law, etc.).

**AI loop**. We refer to *open-loop AI* to models that do not evolve (e.g., inference of a trained supervised model for regression, classification or actuation) or intelligence that does not trigger further analytic (e.g., periodic batch-mode data mining over a data-lake). Open-loop AI may take part in the device (e.g., inference) or the cloud (e.g., transfer learning, federated learning). Open-loop AI techniques are necessary for L2-3 network automation, and will be instrumental also for L4-5.

Conversely, we refer to *closed-loop AI* as to the fact of altering the AI models themselves: we point out that this can happen with any of the *supervised* (e.g., incrementally training a model due to behavior drift of existing classes or appearance of new classes), *unsupervised* (e.g., stream-mode algorithms that alter existing models at any new sample) and *reinforcement* (e.g., continuous exploration phase throughout the whole lifetime of a reinforcement learning, or bandit models) learning classes. Closed-loop AI techniques will be necessary to reach L4-5 network automation.

### C. Network AI Software

Ultimately, the ADN network is executing AI functions which are implemented as software instructions. As any software, AI models need to be designed, maintained and upgraded: even for the simplest L2 O&M task, closed-loop management of AI software is key to successful deployment, which we briefly discuss.
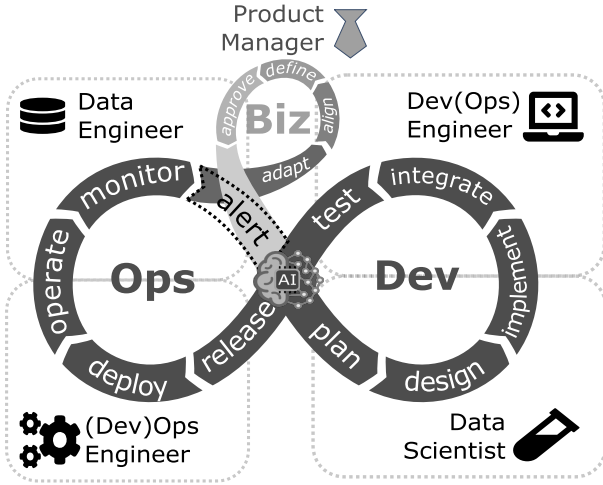
Fig. 3. *Network AI software management*: main actors in the BizDevOps agile loop, with special focus on the AIOps roles and interactions (illustration adapted from [80]).

models may no longer be fit to the environment in which they are deployed, and shall need retraining. In the three-tiered ADN architecture, the offline AI cloud is responsible for tasks such as data lake storage and model catalog management, including model training and maintenance (eg., pruning to fit devices with constrained capabilities, or adaptation in case of unsupported operators on a device, etc.).

## V. CURRENT ACHIEVEMENTS OF NETWORK AI

We now make practical examples of AI-assisted network O&M, that we map over the ADN architecture and illustrate at a glance in Fig.4, ranging from (a) forecast, to (b) traffic management, (c) troubleshooting and (d) auto-tuning. While not exhaustively covering the set of network applications and machine learning techniques[8], the selected set of examples fully span the whole set of supervised, unsupervised, semi-supervised and reinforcement learning branches overviewed in Sec.II, in furthermore open and closed loop modes for both AI and O&M viewpoints.

In this vision paper, we prefer to keep discussion at a qualitative level: we thus avoid embedding quantitative results that are already published elsewhere, to which we rather point the reader to. Additionally, we provide insights from the real problems that AI can find in deployment, that the academic community may not be exposed in their day-to-day work, and is thus less sensitive to: in particular, in each use case we comment AI results under the angles of its (i) generalization capabilities and (ii) benefits vs cost tradeoff – that Sec. III-C and Sec. III-B respectively outlined being of key importance for successful transition from research to products.

### A. Efficiently handling the known (L1 to L2)

Supervised ML techniques, such as regression and classification, are apt at tackling well-specified problems in open-loop O&M settings, to increase visibility about network traffic or distill useful knowledge and information from raw data. In this section, we outline both success and limits for two specific examples of application of each techniques.

**Regression (e.g., QoS/QoE estimation)**. Regression techniques are fit for forecasting, e.g., future traffic demand or user behavior, or for learning complex relationships, such as relating network Quality of Service (QoS) indicator to user Quality of Experience (QoE) as exemplified in Fig. 4-(a). In this latter context, a large body of literature employed ML techniques, to e.g., learn QoS indicators such as latency distribution [82], [83] from topology, traffic matrix and routing information, or learn QoE indicators for specific applications such as Web [84], [85], video [86], [87] or games [88].

AI is desirable in this case, as it can leverage massive volume of data[9] (e.g., automatically collected network [82], [83] or application [85]–[88] QoS/QoE indicators) without

**AIOps Software Lifecycle**. AI software has a peculiar lifecycle, termed AIOps by Gartner, as a particularization of the BizDevOps cycle to take into account specific characteristics of AI-software development. In a nutshell, DevOps is an agile methodology that combines software development (Dev) and IT operations (Ops), to shorten the systems development life cycle by providing continuous delivery, possibly further integrating Business (Biz) aspects. Complementary to the Dev and Ops engineering teams in classic IT and O&M scenarios, AI development requires additional skillsets, which are identified in the Data Engineering and Data Scientist roles respectively, as illustrated in Fig. 3 (further emerging roles are discussed in Sec.VI-A).

**AI Software in the ADN**. Taking the supervised case as an example and with reference to Fig.3, data scientists leverage data gathered by data engineers for model design (e.g., choosing the appropriate ML/AI family, properly train by avoiding overfitting and biases, hyperparameter tuning, etc.). As data science skills may not accessible to all companies, to lower the startup costs, a recent trend is to automate part of the data scientists ML task (e.g., AutoML [81]). Alternatively, the existence of open APIs makes it possible for the emergence of AI marketplaces (e.g., Acumos [79]), that can offer readily trained models, which may lessen the need of data scientists for the most common tasks. With respect to the three-tiered network architecture early illustrated in Fig.2, models gathered by any of the above means can be deployed in edge devices or online flog/cloud (depending on capabilities of the devices).

After models are properly trained and tested, they are deployed by Ops engineer in production system. Data engineers then help monitor data from deployed DL models: tracking of model data is then used to adapt, align and define the priorities (BizDevOps loop) or to simply alert on necessary updates to the model (DevOps loop) for the next Dev phase. As early outlined, due to data drift, or environmental changes, deployed

---

[8]For instance, we plan to address the use of natural language processing techniques for the purpose of network configuration in a separate article.

[9]It may be relevant for readers interested in the Web QoE specific use-case, that we made several datasets, collected in cooperation with Orange Labs [89], [90] and Wikipedia [91] publicly available.
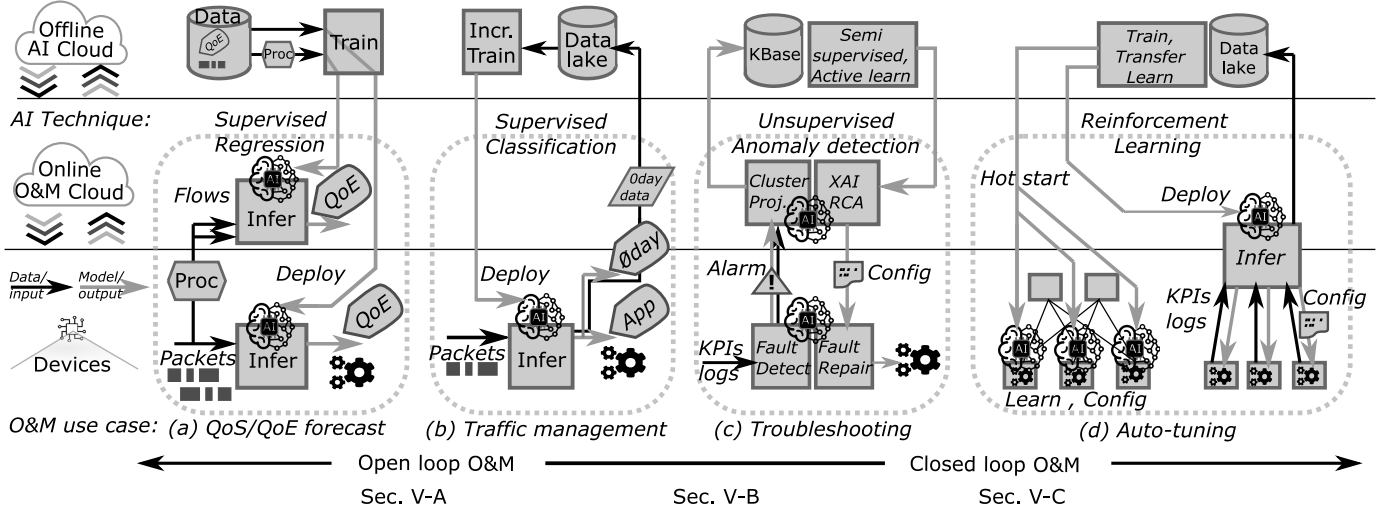
Fig. 4. Synoptic of selected examples in current achievements of AI-assisted networking use cases of: (a) QoS/QoE forecast, (b) traffic management, (c) troubleshooting and (d) auto-tuning.

incurring high labeling cost (unless human opinion is explicitly factored in, as in [84]). On the one hand, work such as [82]–[88] proves that data-driven models provide accurate solutions, at both packet or flow-levels. On the other hand, generalization capabilities and cost bares additional discussion.

As far as generalization is concerned, it is clear that models are tested on a subset of the whole set of applications, games, Webpages, terminals and network conditions. However, real products will be exposed to such diversity: thus, particular attention is needed to stress-test model capabilities beyond the classic techniques (e.g., k-fold cross validation). For instance, as done in [85] for the case of Web QoE, it is useful to systematically analyze model bias induced by training data, and provide guidance for extending data corpus to reduce such via incremental training at higher autonomy levels (L3 and beyond).

Additionally, the academic community often strives to increase accuracy of the proposed solution, irrespective of its computational cost. Yet, as diminishing return has to be expected beyond a given accuracy, solutions that explicitly allow to tradeoff (slight) accuracy loss for (significant) computational savings are to be preferred [85]: the first role of AI researchers is, after all, to judge whether AI is the right tool for the problem at hand – i.e., to avoid that all problems seems nails once you have an (AI) hammer.

**Classification (e.g., Traffic management)**. Traffic management of Fig. 4-(b) is another relevant example where open-loop AI techniques are helpful: traffic prioritization needs coarse-grained traffic category labels, while policing may additional require fine-grained application labels.

In this context, AI is clearly beneficial since encryption is mandating to phase out Deep Packet Inspection (DPI) for behavioral classifiers, with satisfactory accuracy performance often in excess of 95% [92], [93]. Labeled data is in this case well known to be harder to gather and share [94], although the process can be automated to some extent [95]. As the application coverage of publicly available datasets,

is smaller than commercial grade needs, the question about generalization capabilities is manifest: public datasets generally comprise 10-50 classes whereas commercial grade ones [93] include 2000+ labels, with the top-200 classes covering about 95% of the traffic volume. To cope with this mismatch between data accessible to academic vs relevance for industrial needs, we are currently in the process of releasing a highly-anonymized version of our commercial-grade dataset [93] as part of Huawei Rapid Analytics & Model Prototyping (RAMP) data challenges [96].

As for the complexity is concerned, we observe that system researchers [97]–[100] consider extremely simple models (with just 21 neurons [99] or 50 neurons [100] *overall*), whereas AI researchers train excessively big models (state of the art models compared in [92] employ in excess of hundreds-thousands neurons *per-class*). Awareness of commercial-grade challenges and constraints helps landing commercial-grade models out of the lab, by explicitly parsimonious AI-model design (less than hundred thousands neurons for all 200 classes [93]) and optimized implementation (e.g., using domain specific accelerator and languages [101], [102]).

### B. Taming the unknown (L2 to L3)

As Jean Piaget famously said, "Intelligence is not what you know, but what you do when you don't know". With this regard, supervised models are however inherently limited, as they are e.g., unable to guess QoE or labels of completely new applications. Awareness of supervised AI limits is the first step to move up in the autonomy level.

**Anomaly detection (e.g., troubleshooting)**. Troubleshooting of Fig.4-(c) is an example use-case where supervised techniques are not be a good fit. First, as networks strive to operate at very high reliability (i.e., 5-nines), anomalies are rare events (so by definition only very few examples might be available for training). Second, heterogeneity on the collected data and across different networks environments make generalization even harder (and thus unsupervised techniques appealing).

Clearly, the use of lightweight unsupervised algorithms [103], [104] or self-supervised neural network [105], [106] guarantees generalization by design: given the time-varying multi-variate nature of KPI data, both batch-mode [103] (e.g., for periodic analysis) and stream-mode [104], [106] (e.g., for continuous analysis and trigger) algorithms can be leveraged, with stream-mode preferable due to the nature of the application. However, no single anomaly detection algorithm can fully solve the entire end-to-end troubleshooting pipeline [107], which includes further steps such anomaly aggregation across multiple devices (requiring network visibility and thus a better fit for the private O&M cloud). This algorithmic split across device/cloud gives the opportunity to fine-tune algorithmic results with, e.g., semi-supervised XAI [108] or causal RCA [28] approaches, to ameliorate accuracy by exploiting previous information that might be available only for anomalies that occur more often (and that can benefit for global-level knowledge in the cross-network AI cloud).

Along the complexity angle, unless DL models (such as recurrent LSTM [105], VAE [106]) are used for the anomaly detection task, unsupervised techniques are generally lightweight. Yet, it is important to observe that the relative complexity among unsupervised algorithm still spans several orders of magnitude [104], [109], and that devices may be equipped with low computational power (recall Tab.I). Particularly, as telemetry bandwidth is a more stringent bottleneck in this case, algorithms have to consequently be deployed on edge devices for local execution (so that only part of the telemetry is then pushed to online cloud for network-level analysis and visibility), for which saving CPU cycles for similar algorithmic performance is still highly relevant.

**Out-of-distribution detection (e.g., traffic management)**. Supervised technologies remains a good fit for some use-cases, as for instance in the early introduced traffic management of Fig.4-(b): yet, models have a limited knowledge, as they are likely trained with only a fraction of the existing applications. Applications that have never been presented to the model at training are called "out of distribution" (OOD) in AI terms, or "zero-day" (0D) in O&M terms: when presented with OO/0D instances, any supervised model would misclassify them as one of the known classes it has been presented during training. The ability to detect OO/0D samples illustrated in Fig.4-(b) sits as an intermediate step [80], between *open loop* (L1/L2, train once and deploy forever) and *closed loop* operation (L4/L5, to continuously train and deploy models).

By design, complementing supervised models with OO/0D capabilities increase their generality. Both the AI and the O&M communities have come up with general [110] or use-case specific ways [111] to deal with the problem. We point out that for DL architectures, we have contributed a very effective gradient-based technique [93] that does not require architectural changes and works on unmodified models.

At the same time, it is worth reminding that this OO/0D additional feature comes at additional cost: in particular, while the gradient-based technique is faster than feed-forward computation, some of the techniques we experimented with are significantly slower than the DL inference itself [93], and as such are not practical, or need to be used sparingly [80]. Furthermore, whereas OO/0D is useful in solving part of supervised models limits, further effort is needed to assist [112], [113] and explain [114] automated labeling of OO/0D traffic. Additionally, since few OO/0D samples will be initially available, additional techniques (e.g., few-shots [61] and class-incremental [115] learning) will be necessary to close the learning loop.

*C. Learning to learn (L4 and beyond)*

Higher levels of automation imply the use of closed-loop AI techniques in closed-loop O&M settings. As Einstein famously said, "*The true sign of intelligence is not knowledge but imagination*", and to achieve automation at L4 and beyond, the ability to continuously and efficiently learn is key.

**Reinforcement learning (e.g., auto-tuning)**. A classic example of closed-loop O&M is represented by automating resource usage and control, with either centralized fog/cloud agents (as in the WLAN use-case exemplified in the right part of Fig. 4-d), or distributed agents on devices (as in DCN case shown in the left part of Fig.4-d). Note that the choice of distributed vs centralized intelligence may depend on timescale (e.g., DCN) or other consideration (e.g., WLAN Campus AP controller vs individual home APs), thus the examples of Fig.4-(d) are not meant to exhaustively cover all valid possibilities.

In this context, AI agents are employed to reach an objective related to QoS (reducing flow completion time in DCN [63], improving signal quality in WLAN [116], [117]) or QoE (e.g. of videos [118] and games [119]). To attain such goal, agents receive a reward as a result their action (e.g., setting threshold for ECN marking [63], WLAN channel [116] or power [117] configurations, CDN node selection [118], or relative priority of game traffic [119]). In all these disparate cases cases, AI is used to guide the exploration of an otherwise very large state space: e.g., from simpler Stochastic Bandits used in [118], to more complex Deep Reinforcement Learning (DRL) in [63], [116], or Transformers in [117].

Generalization capabilities are important yet hard to ensure, as the environment in which an agent has been trained may differ significantly from the environment where it is deployed. As a rule of thumb, frequent actions give more opportunities to learn: e.g., ECN threshold setting happens on a DCN RTT timescale [63], whereas CDN node selection happens on a per-session basis [118] and WLAN AP configuration on a hourly basis [116]. Additionally, even though algorithms may account for online learning [63], [118], they need to be seeded with an offline training phase [63], [116]. In this offline phase, algorithms are trained with trace-driven approaches [63] or via simulation [116] for hot-start: the more realistic the offline training environment, and the more diverse the environmental conditions explored, the better generalization capabilities can be expected in the real deployment.

In terms of computational complexity, inference is fortunately faster than training, which is expected to be computationally costly. Additionally, in the case of offline training

phase, often the bottleneck is represented by the cost of simulating the environment at each action step (even for DRL [116] and transformer [117] architectures), so that parallel execution is appealing.

**Lifelong learning (e.g., automate model design and update)**. Finally, a complementary viewpoint that encompasses all use-cases illustrated in Fig.4 and is necessary to reach higher automation level is represented by closed-loop AI techniques. In ML what matters is the journey, not the destination: thus successful ADN deployments need to embrace the idea that any model will need to continuously evolve. We illustrate here a number of reasons with the traffic management example of Fig.4-(b) for simplicity.

Techniques under the so called "lifelong learning" [120] umbrella are key for generalization. First, as new zero-day application will keep appear and old applications will be forgotten, there is need for incremental [115] and decremental [121] learning. As existing applications will drift [122], continuous learning will not necessarily only focus on adding new classes, but to update existing ones. As application behavior differ in heterogeneous environment, federated learning [123]–[125] will additionally be needed for privacy or business-sensitive constraints. Lifelong model changes bring significant challenges in terms of learning due to the opposite curses of "catastrophic forgetting" (of previously learned information) vs "intransigence" (to learn new one). Meta-learning can help finding a general representation that assist solving the above problems [120], but as early illustrated, data science skills may not be available at all steps in an organization. As a result, any closed-loop AI technique to automate the data science workflow (e.g., to automatically search for the fittest neural architecture [81] that is easily in/decrementally updated) is relevant in the quest for robust and lifelong generalization of AI models for the ADN.

## VI. FUTURE CHALLENGES FOR NETWORK AI

While the previous section has shown that AI techniques can be profitably used to solve real problems and transferred to real products, however much remains to be done before highly autonomous L5 O&M operations can be attained by the ADN: this section discusses several challenges that need to be solved along the way.

### A. Dirty data

First, it is well accepted that there is no AI without data [56] – so a few points are worth stressing concerning data access, representation and goverance.

**Data access**. While, as earlier introduced, networking data is more fragmented and heterogeneous than natural language or images, it is true that the community-wide effort to share dataset equivalents to image [53] or natural language [54] lags far behind. Challenges [96], [126], [127] are a partial answer to this problem, but the community should recognize the need to federate data collection efforts, as opposite to scatter them along many tiny challenges. Eventually, complementary

to marketplaces for AI models [79], the emergence of data marketplaces [128], clearly abiding to local and international laws as discussed next, is an interesting opportunity to cope with this problem.

**Data representation**. Access to data only solves part of the problem, since is commonly accepted that 50-80% of AI scientist work is spent on data preparation [129], and indeed "lack of data or data quality issues" is the first technical[10] bottleneck to AI adoption, as identified by respondents (5% of which are from the telecommunication and networking industry) to the Oreilly radar survey [64]. Whereas some amount of human assistance to AI models is needed, our goal is not to enslave human labor to reach AI success [59], [60], or at least to consciously use as little human help as possible.

On the one hand, the AI community has come up with a number of best practices and techniques to cope with data quality [130] or lack of labels (via active-learning [62], self-attention [131], few-shot [61] or self-supervision [132]). On the other hand, data in network is intrinsically highly dimensional, topologically complex (several time-varying multi-layer logical graphs), heterogeneous and multi-modal (logs, packets, timeseries, configuration, etc.) and misses a single, unified and universally accepted representation, which remains an open problem.

**Data governance**. Finally, unless data, and meta-data are treasured at (or beyond) the level of first-class citizen for network AI, then much of the AI effort will likely end up being in the 80% of failing projects [2]. The complexity of meta-data management and of properly granting access to data, calls for a more systematic approach to data governance [133]: e.g., by the introduction of data stewards, beyond the roles of data owner, data engineer and data scientist. Data stewards should govern the access to data, on behalf of data owners and using the process developed by data engineers, in compliance with regulation aspects.

### B. AI Regulation

Second, AI will not be successful if it's not legally compliant. As the regulation ecosystem is fragmented as it differ from country to country[11] and additionally evolves over time[12], we briefly review the expected impact on data and AI models.

**Impact on data**. First, we observe that networking O&M area is intrinsecally less sensitive with respect to other industries that treats biometric, physiological or medical data (at the highest risk level for AI Act). Still, even though most content is end-to-end encrypted [138], it is useful reminding that even IP addresses are personal data according to GDPR, and have

---

[10]After the top-3 reasons ("company culture", the "difficulty in identify business use cases", and the "lack of data scientists roles") which are of non-technical nature.

[11]For instance, EU General Data Protection Regulation (GDPR) [134], or the more complex maze of laws across US states [135].

[12]As for the new proposal for a regulatory framework on Artificial Intelligence (AIAct) [136] launched Apr 2021 in EU, or the Personal Information Protection Law (PIPL) [137] introduced in Nov 2021 in China

to be processed accordingly. While some use-cases may raise limited risk, however this is not going to be the general case: as such, lawful compliance of AI data processing should come as a fundamental architectural feature, and not as an afterthought.

For instance, e.g., security protection from intruders may process IP addresses of attackers received by Darknets [139], but may also include part of so called "backscatter" [140] traffic from compromised machines owned by unconscious citizens/companies. As of now, it is safer when data is not shared and stays under same-country (or same-regulation) boundaries: however evolving and heterogeneous regulation between countries can clearly become a headache for AI researchers, and slow down AI adoption for global companies. The early introduced data steward role should additionally be coupled with a legal expert role, clearly mapping data to a specific risk category depending on its processing. Data should also be tainted, so that crossing legal borders should be either lawful or impossible by design in the architecture.

**Impact on models**. In addition to having an impact on data sharing, regulation impacts algorithms even when data is *not* shared: as, for instance, in the case of federated learning [123]. Whereas loads of work ensure privacy of that federated learning process [141], however attacks [142] are still possible that yield to serious privacy leaks – with possible legal consequences, when clients participating in the learning span across multiple countries and regulations.

Moving one step further, the "right to be forgotten" for privacy or business reasons raise the need for *decremental* learning, which has recently enjoyed a number of proposal for Support Vector [121], Random Forest [143] or Neural Networks [144]. However, as pointed out in [145] these algorithms are not immune to attacks, so that machine unlearning can have counterproductive effects on privacy. Further, whereas probabilistic [146] may be enough from a scientist point of view, it would strikes us as odd if this was considered acceptable from a legal standpoint. The lack of clear guidelines and definition can slow-down adoption, for which non-trivial discussion between lawyers and scientists seems in order.

### C. eXplainable AI

Third, AI will not be successful if it's not trusted. In the path to L5 ADN, AI should free humans from the burden of the fast loop (that calls for automation), and facilitating the interaction with them in a slow loop (that calls for explainability). In a sense, AI should smoothly transition from day-to-day technical supervision (i.e., as if AI is a junior colleague that needs coaching) to less and less frequent supervision (i.e., AI becomes a senior respected colleague).

**Fairness and generalization**. Trust can be gained only if AI can provide, in addition to good performance, also clear explanation of the decision process. In other fields where ethical matters are prominent, fairness issues in AI models is well known and debated [147]. Part of the lack of fairness is knowingly rooted in class imbalance in the data used in the first place [148], for which countermeasures exist. If, in the networking field, unfairness is less likely to have life-changing decisions, the existence of such bias can still severely affect the model generalization capabilities. Finally, it has recently been shown that models can be accurate but for the wrong reasons [149], which can rapidly undermine the trust that AI technologies have gained so far. Trust between humans and machines [150], [151] depends on understandability and directability (how easy one can assert, control or influence when something goes wrong): techniques to improve AI explainability [152] are therefore a key element in the future ADN.

**Faithfulness and accountability**. Additionally, whereas most of AI decisions in network O&M will not have dramatic life-changing impacts, such decisions should be accountable from a business and possibly legal (as per previously discussed regulation) perspectives. On the one hand, probabilistic [146] or post-hoc explanation of surrogate models [153], [154] may be sufficient for business-level accountability. On the other hand, proving law compliance at a forensic level may require faithfulness, a key XAI property defined as correctly reflecting the system process for generating the output [155]. This pushes to embed explainability directly in the model, as recent work started advocating [156], [157]. We finally point out that while XAI may come with a cost (e.g., additional computation, or accuracy loss), this seems to be a necessary price to pay in light of the above issues.

### D. Green AI

Fourth, AI revolution will not happen if it's not cost-effective – as we previously observed, efficient hardware accelerators were instrumental to AI success (Sec.III-B), and cost-effective operation is a precondition for business success (Sec.III-D). Interestingly, we observe that the "Green networking" [158] predates by over a decade the corresponding "Green AI" wave [159]. We examine here the complementary aspects of green AI from system and algorithmic perspectives.

**Green system**. As earlier indicated, Gilder law forecasts that bandwidth scaling rate exceeds Moore law. Additionally, OpenAI estimated [160] that the amount of compute used in the largest AI training runs in the last decade has increased exponentially, again faster than Moore law. As such, the carbon impact of AI [161], [162] has rightly come under scrutiny (along with blockhain).

With the example of the traffic management use case early introduced in Sec.V-A, we observe that commercial-grade models can run on DSA (e.g., Ascend310 TPU) supports inference rate at 100Gbps with a power drain of 7W. However, we additionally point out that complementary techniques that avoids AI computation altogether (such as approximate-key-caching [163]) can further bring multiplicative speedups without compromising accuracy. In other technological fields, the awareness of fuel efficiency (for cars) or energy efficiency (for electrical appliances and even light bulbs) is strictly regulated and mandatorily exposed – we believe that the same awareness should be extended to the AI sub-components of

any networked system, by explicitly measuring metrics such as *task/Joule* or *accuracy/Watt*.

**Green models**. Fortunately on the other hand, the AI community has worked toward improving algorithmic efficiency – so that, e.g., the number of floating point operations required to train a classifier to AlexNet-level performance on ImageNet has significantly decreased in the last decade [164].

Always with reference to the traffic management use case of Sec.V-A, recently introduced architectural innovations such as inverted residual as in MobileNet [165] and point-wise convolution as in ShuffleNet [166], significantly reduce model complexity and are amenable to be run on ARM even without any DSA acceleration[13]. At the same time, AI model design is still an art, so that automated ML (autoML) techniques such as Neural Architecture Search (NAS) [81] are appealing to assist AI model design. However, NAS is guided by an indirect measure of computation complexity (i.e., FLOPs). Therefore, further research is needed to explicitly include in the NAS loop more direct metrics, such as speed or energy consumption, to perform an *ecology* of models design space (as well as making the NAS process itself more energy efficient).

### E. AI-Native functions and architecture

Finally, we need to point out that the ADN architecture illustrated and exemplified in this paper is a natural evolution of the current cloud-native network architecture. In other words, exactly as it happened for IP networks, where QoS, mobility, and security features have been added as afterthoughts to the original architecture, the same is now happening for AI technologies. We point out that we are not reviving here the debate of clean-slate vs evolutionary approach [167] to network technology (r)evolution. Neither our goal is to anticipate the success of such changes [168], though we point out that time is ripe in terms of software stack and hardware support, which are known to have a disproportionate impact on technological success [30], [39].

Rather, we believe that this timeframe constitute an opportunity for the network community to rethink the design of an AI-Native architecture, by holistically integrating AI in the whole network landscape, as opposite to just delegate specific tasks to AI islands. In other words, we believe that the future AI-Native architecture should allow to expose, combine and orchestrate explainable AI functions, "wired" to data in a way that respects regulation, overall bringing improvements on network operation at a lower cost and energy footprint.

### ACKNOWLEDGEMENTS

[13]Currently unpublished results show that our ShuffleNet re-implementation on TFLite [49] achieves accuracy on par of [93] by taking $140\mu s$ per classification on a single ARM Cortex A72 core with batch size $b = 1$, and is even faster than on NVIDIA P100 GPU with small batch $b = 8$.

## REFERENCES

[1] N. Feamster *et al.*, "Why (and how) networks should run themselves," 2017.

[2] https://blogs.gartner.com/andrew_white/2019/01/03/our-top-data-and-analytics-predicts-for-2019/.

[3] https://www.oreilly.com/radar/ai-adoption-in-the-enterprise-2020/.

[4] T. Anderson *et al.*, "Overcoming the internet impasse through virtualization," *Computer*, vol. 38, no. 4, pp. 34–41, 2005.

[5] C. Jacquenet, "Optimized, automated, and protective: An operator's view on future networks," *IEEE TNSM*, vol. 18, no. 2, 2021.

[6] *Communication Networks and Service Management in the Era of Artificial Intelligence and Machine Learning*, 2021.

[7] R. Boutaba *et al.*, "A comprehensive survey on machine learning for networking: evolution, applications and research opportunities," *Journal of Internet Services and Applications*, vol. 9, no. 1, 2018.

[8] H. Lutfiyya *et al.*, "SI on Embracing Artificial Intelligence for Network and Service Management," *IEEE TNSM*, vol. 18, no. 4, 2021.

[9] P. Chemouil *et al.*, "SI on Artificial Intelligence and machine learning for networking and communications," *IEEE JSAC*, vol. 37, no. 6, 2019.

[10] ——, "SI on Advances in Artificial Intelligence and Machine Learning for Networking," *IEEE JSAC*, vol. 38, no. 10, 2020.

[11] F. Rosenblatt, "The perceptron. a perceiving and recognizing automation," in *Cornell Aeronautical Laboratory Report 85-460*, 1957.

[12] A. L. Samuel, "Some studies in machine learning using the game of checkers," *IBM Journal of Research and Development*, vol. 3, no. 3, pp. 210–229, 1959.

[13] K. S. Narendra *et al.*, "Application of learning automata to telephone traffic routing and control," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 7, no. 11, pp. 785–792, 1977.

[14] A. Krizhevsky *et al.*, "ImageNet classification with Deep Convolutional Neural Networks," *NIPS*, vol. 25, 2012.

[15] https://twitter.com/matvelloso/status/1065778379612282885.

[16] I. Goodfellow *et al.*, "Generative adversarial nets," *NIPS*, vol. 27, 2014.

[17] https://www.theverge.com/tldr/2018/4/17/17247334/ai-fake-news-video-barack-obama-jordan-peele-buzzfeed.

[18] D. Silver *et al.*, "Mastering the Game of Go without Human Knowledge," *Nature*, vol. 550, no. 7676, pp. 354–359, 2017.

[19] https://openai.com/projects/five/.

[20] T. Mikolov *et al.*, "Distributed representations of words and phrases and their compositionality," in *NIPS*, 2013.

[21] T. B. Brown *et al.*, "Language models are few-shot learners," *NeurIPS*, 2020.

[22] https://openai.com/blog/gpt-3-apps/.

[23] I. Goodfellow *et al.*, *Deep Learning*. MIT Press, 2016, http://www.deeplearningbook.org.

[24] https://en.wikipedia.org/wiki/Netflix_Prize.

[25] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, 2001.

[26] T. Chen *et al.*, "Xgboost: A scalable tree boosting system," in *ACM KDD*, 2016.

[27] F. T. Liu *et al.*, "Isolation forest," in *IEEE ICDM*, 2008.

[28] J. Pearl, "The seven tools of causal inference, with reflections on machine learning," *Commun. of the ACM*, vol. 62, no. 3, 2019.

[29] A. Tavanaei *et al.*, "Deep learning in spiking neural networks," *Neural Networks*, vol. 111, pp. 47–63, 2019.

[30] S. Hooker, "The hardware lottery," *Commun. of the ACM*, vol. 64, no. 12, pp. 58–65, 2021.

[31] J. L. Hennessy *et al.*, "A new golden age for computer architecture," *Commun. of the ACM*, vol. 62, no. 2, pp. 48–60, 2019.

[32] I. Hubara *et al.*, "Binarized neural networks," *NIPS*, vol. 29, 2016.

[33] M. Rastegari *et al.*, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *ECCV*, 2016, pp. 525–542.

[34] M. Davies, "Lessons from Loihi: Progress in Neuromorphic Computing," in *IEEE Symposium on VLSI Circuits*, 2021.

[35] https://www.tensorflow.org/.

[36] https://github.com/mindspore-ai/mindspore.

[37] https://pytorch.org/.

[38] https://scikit-learn.org/.

[39] P. Molino *et al.*, "Declarative machine learning systems," *Commun. of the ACM*, vol. 65, no. 1, pp. 42–49, 2022.

[40] https://developer.arm.com/ip-products/processors/cortex-a/cortex-a72.

[41] https://coral.ai/docs/edgetpu/benchmarks/.

[42] https://e.huawei.com/en/products/cloud-computing-dc/atlas/ascend-310.

[43] https://images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf.

[44] https://cloud.google.com/tpu.

[45] https://www.hisilicon.com/en/products/Ascend/Ascend-910.

[46] https://www.nvidia.com/en-us/data-center/products/egx-converged-accelerator/.

[47] https://github.com/ARM-software/armnn.

[48] L. Linguaglossa *et al.*, "Survey of performance acceleration techniques for network function virtualization," *Proceedings of the IEEE*, vol. 107, no. 4, pp. 746–764, Apr. 2019.

[49] https://www.tensorflow.org/lite.

[50] R. David *et al.*, "Tensorflow lite micro: Embedded machine learning for tinyml systems," *MLSys*, vol. 3, 2021.

[51] http://yann.lecun.com/exdb/mnist/.

[52] https://www.cs.toronto.edu/~kriz/cifar.html.

[53] https://www.image-net.org/.

[54] https://commoncrawl.org/.

[55] J. Jumper *et al.*, "Highly accurate protein structure prediction with alphafold," *Nature*, vol. 596, no. 7873, pp. 583–589, 2021.

[56] C. Gröger, "There is no AI without data," *Commun. of the ACM*, Nov 2021.

[57] https://www6.slac.stanford.edu/research/scientific-computing.

[58] https://home.cern/resources/faqs/facts-and-figures-about-lhc.

[59] https://blog.cloudflare.com/introducing-cryptographic-attestation-of-personhood/.

[60] https://anatomyof.ai/index.html.

[61] Y. Wang *et al.*, "Generalizing from a few examples: A survey on few-shot learning," *ACM Computing Surveys*, vol. 53, no. 3, pp. 1–34, 2020.

[62] B. Settles, "Active learning literature survey," 2009.

[63] S. Yan *et al.*, "ACC: Automatic ECN tuning for high-speed datacenter networks," in *ACM SIGCOMM*, 2021.

[64] https://www.oreilly.com/radar/ai-adoption-in-the-enterprise-2020/.

[65] https://www.oreilly.com/radar/ai-adoption-in-the-enterprise-2021/.

[66] https://www.etsi.org/technologies/experiential-networked-intelligence.

[67] https://www.etsi.org/technologies/zero-touch-network-service-management.

[68] https://datatracker.ietf.org/rg/coinrg/about/.

[69] https://e.huawei.com/en/products/enterprise-networking/routers/ne/ne8000.

[70] https://e.huawei.com/en/products/enterprise-networking/wlan/wifi-6/.

[71] https://e.huawei.com/en/products/enterprise-networking/security/firewall-gateway/usg6500e.

[72] https://e.huawei.com/cn/products/network-management-and-analysis-software/imaster-nce-campus.

[73] https://carrier.huawei.com/cn/products/wireless-network/mae.

[74] https://www.hwtelcloud.com/.

[75] C. Lao *et al.*, "ATP: In-network aggregation for multi-tenant learning." in *USENIX NSDI*, 2021.

[76] J. Fei *et al.*, "Efficient sparse collective communication and its application to accelerate distributed deep learning," in *ACM SIGCOMM*, 2021.

[77] https://support.huaweicloud.com/modelarts/index.html.

[78] https://aws.amazon.com/sagemaker/.

[79] https://www.acumos.org/.

[80] L. Yang *et al.*, "Quality monitoring and assessment of deployed Deep Learning models for network AIOps," *IEEE Network*, 2021.

[81] T. Elsken *et al.*, "Neural architecture search: A survey," *The Journal of Machine Learning Research*, vol. 20, no. 1, pp. 1997–2017, 2019.

[82] S. Xiao *et al.*, "Deep-q: Traffic-driven qos inference using deep generative network," in *ACM SIGCOMM, NetAI Workshop*, 2018.

[83] K. Rusek *et al.*, "Unveiling the potential of graph neural networks for network modeling and optimization in SDN," in *ACM Symposium on SDN Research (SOSR)*, 2019.

[84] F. Salutari *et al.*, "Analyzing wikipedia users' perceived quality of experience: A large-scale study," *IEEE TSNM*, vol. 17, no. 2, 2020.

[85] A. Huet *et al.*, "Deployable models for approximating web qoe metrics from encrypted traffic," *IEEE TNSM*, 2021.

[86] D. Ghadiyaram *et al.*, "Learning a continuous-time streaming video qoe model," *IEEE Trans. on Image Processing*, vol. 27, no. 5, 2018.

[87] T. Mangla *et al.*, "eMIMIC: Estimating HTTP-Based Video QoE Metrics from Encrypted Network Traffic," in *IFIP TMA*, 2018.

[88] S. Zadtootaghaj *et al.*, "NDNetGaming – development of a no-reference deep CNN for gaming video quality prediction," *Multimedia Tools and Applications*, 2020.

[89] https://figshare.com/articles/dataset/Revealing_QoE_of_Web_Users_from_Encrypted_Network_Traffic/12459293.

[90] https://figshare.com/articles/dataset/Detecting_Degradation_of_Web_Browsing_Quality_of_Experience/13089854.

[91] https://figshare.com/articles/dataset/A_Large_Scale_Study_of_Wikipedia_Users_Perceived_Quality_of_Experience/11365607.

[92] G. Aceto *et al.*, "Mobile encrypted traffic classification using deep learning," in *IEEE TMA*, 2018.

[93] L. Yang *et al.*, "Deep learning and traffic classification: Lessons learned from a commercial-grade dataset with hundreds of encrypted and zero-day applications," *IEEE TNSM*, vol. 18, 2021.

[94] A. W. Moore *et al.*, "Toward the accurate identification of network applications," in *PAM*, 2005.

[95] G. Aceto *et al.*, "Mirage: Mobile-app traffic capture and ground-truth creation," in *IEEE ICCCS*, 2019.

[96] http://xianti.fr.

[97] Z. Xiong *et al.*, "Do switches dream of machine learning?: Toward in-network classification," in *ACM HotNets XVIII*, 2019.

[98] C. Busse-Grawitz *et al.*, "pForest: In-Network Inference with Random Forests," *arXiv:1909.05680*.

[99] T. Swamy *et al.*, "Taurus: An intelligent data plane," *arXiv:2002.08987*, 2020.

[100] G. Siracusano *et al.*, "Running neural networks on the NIC," *arXiv:2009.02353*, 2020.

[101] M. Gallo *et al.*, "Fenxi: Deep-learning traffic analytics at the edge," *ACM/IEEE Symposium on Edge Computing (SEC)*, 2021.

[102] ——, "Real-time deep learning based traffic analytics," in *ACM SIGCOMM, Demo session*, Aug. 2020.

[103] A. Putina *et al.*, "Random histogram forest for unsupervised anomaly detection," in *IEEE ICDM*, 2020.

[104] ——, "Online anomaly detection leveraging stream-based clustering and real-time telemetry," *IEEE TNSM*, vol. 18, no. 1, 2021.

[105] P. Malhotra *et al.*, "LSTM-based encoder-decoder for multi-sensor anomaly detection," *ICML*, 2016.

[106] H. Xu *et al.*, "Unsupervised anomaly detection via variational auto-encoder for seasonal KPIs in web applications," in *WWW*, 2018.

[107] T. Li *et al.*, "Flap: An end-to-end event log analysis platform for system management," in *ACM KDD*, 2017.

[108] J. M. Navarro *et al.*, "Hurra! human readable router anomaly detection," Sep. 2020.

[109] ——, "Human readable network troubleshooting based on anomaly detection and feature scoring," Tech. Rep., arXiv:2108.11807, Mar 2021.

[110] M. A. Pimentel *et al.*, "A review of novelty detection," *Signal Processing*, vol. 99, pp. 215–249, 2014.

[111] J. Zhang *et al.*, "Robust network traffic classification," *IEEE/ACM TNET*, vol. 23, no. 4, pp. 1257–1270, jul 2015.

[112] ——, "Autonomous unknown-application filtering and labeling for dl-based traffic classifier update," in *IEEE INFOCOM*, 2020.

[113] T. van Ede *et al.*, "FlowPrint: Semi-supervised mobile-app fingerprinting on encrypted network traffic," in *Network and Distributed System Security Symposium (NDSS)*, 2020.

[114] C. Beliard *et al.*, "Opening the deep pandora box: Explainable traffic classification," in *IEEE INFOCOM, Demo session*, july 2020.

[115] G. Bovenzi *et al.*, "A first look at class incremental learning in deep learning mobile traffic," in *IFIP TMA*, 2021.

[116] O. Iacoboaiea *et al.*, "Real-Time channel management in WLANs: DRL versus heuristics," in *IFIP Networking*, 2021.

[117] M. Boffa *et al.*, "Neural combinatorial optimization beyond the tsp: Existing architectures under-represent graph structure," in *AAAI, GLCR workshop*, 2022.

[118] J. Jiang *et al.*, "Pytheas: Enabling data-driven quality of experience optimization using group-based exploration-exploitation," in *USENIX NDSI*, 2017.

[119] G. Sviridov *et al.*, "Removing human players from the loop: Ai-assisted assessment of gaming qoe," in *IEEE INFOCOM, NI Workshop*, 2020.

[120] Z. Chen *et al.*, "Lifelong machine learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 12, no. 3, 2018.

[121] G. Cauwenberghs *et al.*, "Incremental and decremental support vector machine learning," *NeurIPS*, pp. 409–415, 2001.

[122] V. Carela-Español *et al.*, "A streaming flow-based technique for traffic classification applied to 12+ 1 years of internet traffic," *Telecommunication Systems*, vol. 63, no. 2, pp. 191–243, 2016.

[123] J. Konečnỳ *et al.*, "Federated learning: Strategies for improving communication efficiency," in *NeurIPS Workshop on Private Multi-Party Machine Learning*, 2016.

[124] L. Yang *et al.*, "Heterogeneous data-aware federated learning," in *IJCAI Workshop on Federated Learning*, 2020.

[125] https://www.hwtelcloud.com/products/fed.

[126] https://www.itu.int/en/ITU-T/AI/challenge/2020/.

[127] https://challenge.aiforgood.itu.int/.

[128] X.-Y. Li *et al.*, "Can china lead the development of data trading and sharing markets?" *Commun. of the ACM*, vol. 61, no. 11, 2018.

[129] https://www.nytimes.com/2014/08/18/technology/for-big-data-scientists-hurdle-to-insights-is-janitor-work.html.

[130] W. Kim *et al.*, "A taxonomy of dirty data," *Data mining and knowledge discovery*, vol. 7, no. 1, pp. 81–99, 2003.

[131] A. Vaswani *et al.*, "Attention is all you need," in *NIPS*, 2017.

[132] C. Doersch *et al.*, "Multi-task self-supervised visual learning," in *IEEE ICCV*, 2017, pp. 2051–2060.

[133] R. Abraham *et al.*, "Data governance: A conceptual framework, structured review, and research agenda," *International Journal of Information Management*, vol. 49, pp. 424–438, 2019.

[134] https://en.wikipedia.org/wiki/General_Data_Protection_Regulation.

[135] https://www.endpointprotector.com/blog/eu-vs-us-how-do-their-data-protection-regulations-square-off/.

[136] https://www.europarl.europa.eu/thinktank/en/document/EPRS_BRI(2021)698792.

[137] https://en.wikipedia.org/wiki/Personal_Information_Protection_Law_of_the_People's_Republic_of_China.

[138] A. P. Felt *et al.*, "Measuring HTTPS adoption on the Web," in *USENIX Security Symposium*, 2017, pp. 1323–1338.

[139] L. Gioacchini *et al.*, "Darkvec: Automatic analysis of darknet traffic with word embeddings," in *ACM CoNEXT*, 2021.

[140] E. Balkanli *et al.*, "On the analysis of backscatter traffic," in *IEEE LCN*, 2014.

[141] S. Truex *et al.*, "A hybrid approach to privacy-preserving federated learning," in *ACM Workshop on AI and Security*, 2019.

[142] Z. Wang *et al.*, "Beyond inferring class representatives: User-level privacy leakage from federated learning," in *IEEE INFOCOM*, 2019.

[143] J. Brophy *et al.*, "Machine unlearning for random forests," in *ICML*, 2021.

[144] S. Neel *et al.*, "Descent-to-delete: Gradient-based methods for machine unlearning," pp. 931–962.

[145] M. Chen *et al.*, "When machine unlearning jeopardizes privacy," in *ACM CCS*, 2021.

[146] D. M. Sommer *et al.*, "Towards probabilistic verification of machine unlearning," 2022.

[147] https://www.nytimes.com/2020/12/03/technology/google-researcher-timnit-gebru.html.

[148] N. Mehrabi *et al.*, "A survey on bias and fairness in machine learning," *ACM Computing Surveys*, vol. 54, no. 6, pp. 1–35, 2021.

[149] S. Lapuschkin *et al.*, "Unmasking Clever Hans predictors and assessing what machines really learn," *Nature*, vol. 10, no. 1, 2019.

[150] R. R. Hoffman *et al.*, "Trust in automation," *IEEE Intelligent Systems*, vol. 28, no. 1, pp. 84–88, 2013.

[151] N. Wang *et al.*, "Trust calibration within a human-robot team: Comparing automatically generated explanations," in *ACM/IEEE HRI*, 2016.

[152] A. B. Arrieta *et al.*, "Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI," *Information Fusion*, vol. 58, pp. 82–115, 2020.

[153] Z. Meng *et al.*, "Interpreting deep learning-based networking systems," in *ACM SIGCOMM*, 2020.

[154] S. M. Lundberg *et al.*, "A unified approach to interpreting model predictions," in *NIPS*, 2017.

[155] P. J. Phillips *et al.*, "Four principles of explainable artificial intelligence," *US National Institute of Standards and Technology (NIST) Draft*, 2020.

[156] D. Alvarez-Melis *et al.*, "Towards robust interpretability with self-explaining neural networks," *NeurIPS*, 2018.

[157] C. Rudin, "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead," *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206–215, 2019.

[158] A. P. Bianzino *et al.*, "A survey of green networking research," *IEEE Commun. Surveys & Tutorials*, vol. 14, no. 1, pp. 3–20, 2010.

[159] R. Schwartz *et al.*, "Green AI," *Commun. of the ACM*, vol. 63, no. 12, pp. 54–63, 2020.

[160] https://openai.com/blog/ai-and-compute/.

[161] P. Dhar, "The carbon impact of artificial intelligence," *Nature Machine Intelligence*, vol. 2, no. 8, pp. 423–425, 2020.

[162] https://www.wired.com/story/ai-great-things-burn-planet/.

[163] A. Finamore *et al.*, "Accelerating deep learning classification with error-controlled approximate-key caching," *IEEE INFOCOM*, 2022.

[164] https://openai.com/blog/ai-and-efficiency.

[165] A. Howard *et al.*, "Inverted residuals and linear bottlenecks: Mobile networks for classification, detection and segmentation," in *IEEE CVPR*, 2018.

[166] N. Ma *et al.*, "Shufflenet v2: Practical guidelines for efficient cnn architecture design," in *ECCV*, 2018.

[167] J. Rexford *et al.*, "Future internet architecture: clean-slate versus evolutionary research," *Commun. of the ACM*, vol. 53, no. 9, 2010.

[168] M. Ammar, "ex uno pluria: The service-infrastructure cycle, ossification, and the fragmentation of the internet," *ACM SIGCOMM Computer Communication Review*, vol. 48, no. 1, 2018.

**Dario Rossi** is Director of Huawei AI4NET Lab and Director of the DataCom Department at the Paris Research Center, France. Before joining Huawei in 2018, he held Full Professor positions at Telecom Paris and Ecole Polytechnique and was holder of Cisco's Chair NewNet Paris. He has coauthored 15 patents and over 200 papers in leading conferences and journals, that received 9 best paper awards, a Google Faculty Research Award (2015) and an IRTF Applied Network Research Prize (2016). He is a Senior Member of IEEE and ACM.

**Liang Zhang** is Vice-Director of Huawei AI4NET Lab and Director of the DataCom AI Department at the Nanjing Research Center, China. He received the PhD degree from Southeast University, Nanjing, China, in 2010. His research interests include intelligent fault analysis, network traffic analysis and network optimization.