

Fog and IoT: An Overview of Research Opportunities

Mung Chiang, *Fellow, IEEE*, and Tao Zhang, *Fellow, IEEE*

(Invited Paper)

Abstract—Fog is an emergent architecture for computing, storage, control, and networking that distributes these services closer to end users along the cloud-to-things continuum. It covers both mobile and wireline scenarios, traverses across hardware and software, resides on network edge but also over access networks and among end users, and includes both data plane and control plane. As an architecture, it supports a growing variety of applications, including those in the Internet of Things (IoT), fifth-generation (5G) wireless systems, and embedded artificial intelligence (AI). This survey paper summarizes the opportunities and challenges of fog, focusing primarily in the networking context of IoT.

Index Terms—Edge computing, edge networking, edge storage, fog, fog computing, fog control, fog networking, fog storage, Internet of Things (IoT).

I. INTRODUCTION

FOG IS an architecture that distributes computation, communication, control and storage closer to the end users along the cloud-to-things continuum. Sometimes the term “fog” is used interchangeably with the term “edge,” although fog is broader than the typical notion of edge. The relevance of fog/edge is rooted in both the inadequacy of the traditional cloud and the emergence of new opportunities for the Internet of Things, 5G and embedded artificial intelligence.

Over the past decade, moving computing, control, and data storage into the cloud has been an important trend. In particular, computing, storage, and network management functions are shifted to centralized data centers, backbone IP networks, and cellular core networks. Today, however, cloud computing is encountering growing challenges in meeting many new requirements in the Internet of Things (IoT).

At the same time, there has also been a surging number and variety of powerful end-user, network edge, and access devices: smartphones, tablets, smart home appliances, small cellular base stations, edge routers, traffic control cabinets along the roadside, connected vehicles, smart meters, and energy controllers in a smart power grid, smart building controllers, industrial control systems, just to name a few.

Manuscript received April 29, 2016; revised May 30, 2016; accepted June 2, 2016. Date of publication June 23, 2016; date of current version January 10, 2017.

M. Chiang is with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: chiangm@princeton.edu).

T. Zhang is with Cisco Systems, San Jose, CA 95134 USA (e-mail: tazhang2@cisco.com).

Digital Object Identifier 10.1109/IIOT.2016.2584538

Many more smart clients and edge devices, such as drones, industrial and consumer robots, information-transmitting light-bulbs, computers on a stick, and button-sized radio frequency tuners, are following right behind.

It has therefore become feasible and interesting to ask: “What can be done closer to the end users?” Can your car become your primary data store? Can a single appliance in your house integrate the different services and applications that have been provided by separate systems such as TV set-boxes, home media centers, Internet access routers, and smart energy control boxes? What if smartphones themselves can collectively perform radio network control functions that are performed by gateways in the LTE core networks today? What can a crowd of nearby smart endpoints and network edge devices collectively accomplish through a distributed and immersive network on the edge? Can smart edge devices collectively enable ultralow or even deterministic latency to support delay-sensitive applications, such as real-time data analytics on the edge, mining of streaming data, and industrial control functions?

What these questions point to is a pendulum swinging now back from “click” toward “brick,” from “more centralization” to “more immersive distribution,” from “bigger and farther away” clouds to not just smaller clouds but computation and control closer to sensors, actuators, and users. The pendulum between centralization and distribution is decades-old, with two distinct flavors of “distribution”: first is the end-to-end principle as exemplified by TCP congestion control and perhaps peer-to-peer (P2P) multicast overlay, and second is leveraging local proximity as in Ethernet and sensor networks. Fog embodies and further accelerates this click-to-brick swing-back from both angles, and for not only the data plane but also the control plane.

This paper starts with the range of new challenges in the emerging IoT and the difficulty to address these challenges with today’s computing and networking models. It then discusses why we will need a new architecture—fog for computing, storage, networking, and control—and how it can fill the technology gaps and create new business opportunities.

Architecture is about functionality allocation [1]: deciding who does what and how to “glue” them back together. Unlike the more mature technology fields such as serial computation, digital communication, and the Internet, where strong and solid architectural foundation has been laid, we are still searching for architectural principles for many emerging systems and applications such as IoT, cyber-physical systems, and

embedded artificial intelligence (AI). We need to make fundamental decisions ranging from where to compute and where to store data along the “cloud-to-things” continuum to how to map computation tasks into a substrate of heterogeneously capable and variably available nodes. Fog provides a direction for us to explore such an architecture; and this paper pays particular attention to IoT as a large application domain over the fog architectural foundation.

II. NEW CHALLENGES IN IoT REQUIRES NEW ARCHITECTURE

The emerging IoT introduces many new challenges that cannot be adequately addressed by today’s cloud and host computing models alone. Here, we discuss several such fundamental challenges.

A. Stringent Latency Requirements

Many industrial control systems, such as manufacturing systems, smart grids, oil and gas systems, and goods packaging systems, often demand that end-to-end latencies between the sensor and the control node stay within a few milliseconds [11]. Many other IoT applications, such as vehicle-to-vehicle communications, vehicle-to-roadside communications, drone flight control applications, virtual reality applications, gaming applications, and real-time financial trading applications, may require latencies below a few tens of milliseconds. These requirements fall far outside what mainstream cloud services can achieve.

B. Network Bandwidth Constraints

The vast and rapidly growing number of connected things is creating data at an exponential rate [12]. A connected car, for example, can create tens of megabytes of data per second. This will include data about: 1) the car’s mobility such as its routes and speeds; 2) the car’s operating conditions such as the wear and tear on its components; 3) the car’s surrounding environment such as road and weather conditions; and 4) videos recorded by the car’s safety cameras. An autonomous vehicle will generate even more data, which was estimated to be about one gigabyte per second [13]. The U.S. smart grid is expected to generate 1000 petabytes of data each year. By comparison, the U.S. Library of Congress generated about 2.4 petabytes of data a month, Google trafficked about one petabyte a month, and AT&T’s network consumed 200 petabytes a year in 2010 [14].

Sending all the data to the cloud will require prohibitively high network bandwidth. It is often unnecessary or sometimes prohibited due to regulations and data privacy concerns. ABI Research estimates that 90% of the data generated by the endpoints will be stored and processed locally rather than in the cloud [12].

C. Resource-Constrained Devices

Many IoT devices will have severely limited resources. Examples include sensors, data collectors, actuators, controllers, surveillance cameras, cars, trains, drones, and medical devices embedded in patients.

Many resource-constrained devices will not be able to rely solely on their own limited resources to fulfill all their computing needs. Requiring all of them to interact directly with the cloud will be unrealistic and cost prohibitive as well, because such interactions often require resource-intensive processing and complex protocols. For example, the multitude of micro-computers on a modern vehicle need firmware updates, but requiring each of these resource-constrained devices to perform the heavy cryptographic operations and sophisticated procedures required to obtain firmware updates from cloud services will be impractical.

D. Cyber-Physical Systems

As more cyber-physical systems are connected to the IoT, the pendulum between the brick versus the click is starting to swing back toward the brick again, where interactions, and often times close integrations, between cyber systems and physical systems are becoming increasingly important and bring new business priorities and operational requirements. Examples of cyber-physical systems include industrial control systems, smart cities, and connected cars and trains. In such systems, uninterrupted and safe operation is often the top priority. Taking a system offline for any reason can cause significant business loss or intolerable customer inconvenience, and therefore, must be planned days, weeks, and even months in advance in some cases [18].

- 1) Requiring cars to be brought to repair shops just to install software update packages can cause intolerable inconvenience and result in heavy cost to both car owners and carmakers.
- 2) A nuclear reactor typically runs on 18-month cycles and any downtime can cause tens of thousands of dollars [16].
- 3) Many other industrial control or manufacturing systems, such as car assembly plants and electrical power generators in the energy grids, have similar requirements for uninterrupted safe operations and require weeks to months lead times to plan for system down times.

As a result, unlike the routers, switches, personal computers, and smartphones in today’s Internet, the timings and opportunities for updating the hardware and software in such cyber-physical systems can be severely limited. Many time-critical control applications, which need to be updated over time, cannot be moved to the cloud due to delay, bandwidth, or other constraints. Therefore, a new computing and networking architecture will be needed to reduce the needs for the hardware and software in mission-critical systems to be updated over time.

E. Uninterrupted Services With Intermittent Connectivity to the Cloud

Cloud services will have difficulty providing uninterrupted services to devices and systems that have intermittent network connectivity to the cloud. Such devices include vehicles, drones, and oil rigs. For example, an oil rig in the ocean and far away from shore may have only satellite communication channels to connect to the cloud. These satellite channels can

suffer widely fluctuating quality and intermittent availability. However, applications such as data collection, data analytics, and controls for the oil rig have to be available even when the rig does not have network connectivity with the cloud. As another example, when a car traverses an area where it loses Internet connectivity, many services and applications for the devices and people in the car must continue to be available. When a car breaks down in such an area and needs to have one of its electronic control unit (ECU) replaced before it can run again, the new ECU should be authenticated to prevent any unauthorized and potentially malware-infected ECUs from being installed on the vehicle. However, cloud-based authentication services will not be available in this scenario.

F. New Security Challenges

Existing cyber security solutions for today's Internet, designed primarily for protecting enterprise networks, data centers, and consumer electronics, have focused on providing perimeter-based protections. In particular, a system or an individual device under protection is placed behind firewalls that work with intrusion detection and prevention systems to prevent security threats from breaking through the protected perimeters. Some resource-intensive security functions are also being moved to the cloud. Existing cloud-based security services continue to focus on providing perimeter-based protection, such as redirecting email and Web traffic to the clouds for threat detection, and redirecting access control requests to the clouds for authentication and authorization processing. Should threats penetrate these protections, the common responses have been for human operators to take the system offline, clean up or replace compromised files and devices, and then put the system back online.

This existing security paradigm will no longer be adequate for addressing many new security challenges in the emerging IoT. Here, we discuss several such challenges.

1) *Keeping Security Credentials and Software up to Date on Large Number of Devices:* As the number and variety of the connected devices increase, a growing challenge will be how to manage the security credentials on these devices and how to keep the security credentials and security software on the devices up to date. Requiring every device to connect to the cloud to update its security credentials and software will be impractical.

2) *Protecting Resource-Constrained Devices:* Many resource-constrained devices in the IoT will not have sufficient resources to protect themselves adequately. These devices may have very long lifespans, and the hardware and software on them can be impractical to upgrade. Yet, these devices will need to remain secure over their long lifespans. For example, replacing any hardware on cars, which have already been sold to consumers, can create significant inconvenience to vehicle owners and result in heavy costs and reputation damages to carmakers. However, over a car's long lifespan that averages about 11.4 years [17], security threats will become significantly more advanced, many new threats will appear, and the mechanisms required to combat the growing threats will need to be enhanced and upgraded accordingly. Therefore, a fundamental question arises: How to

protect a very large number of resource-constrained devices from security attacks?

3) *Assessing the Security Status of Large Distributed Systems in Trustworthy Manner:* IoT will support many large distributed systems. A connected transportation system, for example, may have thousands of devices deployed throughout a city to control traffic signals and communicate with vehicles. A large carmaker will need to ensure the security of tens of millions of cars on the road in a large country such as the USA. An oil and gas company may need to interconnect hundreds of remote sites such as oil rigs, exploration sites, refineries, and pipelines. A smart grid will consist of networked subsystems for metering, data collection, data aggregation, energy distribution, and demand response in multiple geographical areas.

Therefore, the ability to tell, in a trustworthy manner, whether a large number of distributed devices and systems are operating securely, will be essential. However, conventional approaches have difficulty meeting both the scalability and the trustworthy monitoring requirements at the same time.

Today's security health monitoring systems rely on collecting security status messages and log data from devices. These systems, however, can often generate untrustworthy results when applied in some IoT systems.

- 1) Many devices operating in physically unprotected environments can be compromised and used to send false information [22]–[24]. Adversaries can also easily use these compromised devices to form a local majority in many IoT scenarios. For example, they may compromise the majority of the smart meters in a house, a building, or even an entire region. As a result, existing mechanisms for detecting false information, which typically rely on the majority of the data sources to be honest (i.e., uncompromised and not malfunctioning), will no longer be adequate.
- 2) Attackers can compromise a cyber-physical system and damage the physical equipment while keeping the messages to and from the system appear normal. A prime example is the Stuxnet attack on the Iranian nuclear facility—the Stuxnet worm masqueraded the attack by sending normal status messages to the system administrators while spinning the nuclear reactor out of control [19]–[21].

To increase the trustworthiness of security status monitoring, remote attestation mechanisms allow a device to cryptographically prove its trustworthiness to a remote verifier [25], [26]. A device makes a claim about certain properties of its hardware, software, or runtime environment to the verifier and uses its security credentials (e.g., a hardware-based root of trust and public key certificates) to vouch for these properties. The verifier then cryptographically verifies these claims.

However, existing remote attestation methods have focused on enabling an individual device to attest to its own trustworthiness. Many resource-constrained devices in the IoT will not be able to support processing-intensive remote attestation. Even when they can, forcing a large number of devices to perform remote attestation can result in prohibitively

high cost and management complexity. Furthermore, existing remote attestation technology alone cannot handle the case where a device itself is not compromised but its sensory input is.

4) *Responding to Security Compromises Without Causing Intolerable Disruptions:* Today's incident response solutions rely predominately on brute-force mechanisms such as shutting down a potentially compromised system, reinstalling and rebooting its software, or replacing its components and subsystems. Such highly disruptive responses, which largely disregard how severe the compromises actually are, can cause intolerable disruptions to mission-critical systems. However, maintaining uninterrupted and safe operation, even when the system is compromised, is often the highest priority for mission-critical systems such as industrial control systems, manufacturing plants, connected vehicles, drones, and smart grids.

- 1) An electric power generator may be infected by a malware that merely seeks to steal power for unauthorized use. Shutting down the power generator could cause severe disruptions to the smart grid and excessive power outages.
- 2) Industrial control systems often have little tolerance for down time. Manufacturing operations can also have critical safety implications. As a result, manufacturers usually value uninterrupted operation and safety over system integrity. This means that hardware and software updates can only be installed during a system's scheduled down times, which have to be short and far between, rather than every time any security compromise is detected.
- 3) A connected car can be infected by malware that can become active while the car is in motion. While the malware can do a range of damages to the vehicle and can put the driver and passengers in harm's way, abruptly shutting down the engine each time any malware is detected could be an even quicker and surer way to cause deadly traffic accidents.
- 4) If a drone flying midair is abruptly turned off just because a security compromise is detected, it can crash from the sky onto people, houses, and other properties to cause serious damages. Instead, safe landing or safe return-home mechanisms will be essential for responding to such security threats that can compromise a drone's flight.
- 5) A server in a data center may be infected by a spyware that seeks to steal commercial secrets. While allowing such a compromised server to continue to operate could give the attacker access to some sensitive data, it may not directly impact the data-center's mission-critical services. If we shut down the server, or halt the execution of the malware-infected files to wait for the malware to be removed, the system downtime could cause significantly more damage, including causing vast economic losses to the data center operator, business disruptions to those who count on the data centers to operate their businesses, and inconvenience to other users of the data center.

Therefore, today's highly disruptive incident response paradigm will no longer be adequate for securing the many mission-critical systems in the emerging IoT.

III. AN EMERGING ERA OF FOG

Filling the technology gaps in supporting IoT will require a new architecture—fog—that distributes computing, control, storage, and networking functions closer to end user devices.

Complementing the centralized cloud, fog stands out along the following three dimensions:

- 1) Carry out a substantial amount of data storage at or near the end user (rather than storing data only in remote data centers).
- 2) Carry out a substantial amount of computing and control functions at or near the end user (rather than performing all these functions in remote data centers and cellular core networks). Such computing and control functions can include the following.
 - a) Applications for end users and their devices.
 - b) Functions for controlling and operating end-user systems such as manufacturing systems, vehicles, and smart grids.
 - c) Services for managing end-user as well as end-to-end networks, systems, and applications.
 - d) Services for supporting cloud-based applications, such as collecting and preprocessing data to be sent to the cloud.
- 3) Carry out a substantial amount of communication and networking at or near the end user (rather than routing all network traffic through the backbone networks). This can include, for example, ways to improve the performance and scalability of local D2D networks, intelligent control of radio access networks (RANs), organize and manage local mobile ad-hoc networks, and integrate local ad-hoc networks with the infrastructure networks.

Fog and cloud complement each other to form a service continuum between the cloud and the endpoints by providing mutually beneficial and interdependent services to make computing, storage, control, and communication possible anywhere along the continuum.

- 1) *Fog Enables a Service Continuum:* Fog fills the gap between the cloud and the things to enable a service continuum. For example, to the wearable devices, a mobile phone may become the fog to provide local control and analytics applications to the wearable devices. When the user is inside her vehicle, the vehicle can become the fog for her mobile phone to allow many smartphone functions, such as display, user interface, audio, phone book, to be moved to the vehicle. Roadside traffic control equipment can in turn serve as the fog for the vehicle to provide traffic information to the vehicle.
- 2) *Fog and Cloud Are Interdependent:* For example, cloud services may be used to manage the fog. Fog can act as the proxy of the cloud to deliver cloud services to endpoints, and act as the proxy of the endpoints to interact with the cloud. Furthermore, fog can be the beachheads for collecting and aggregating data for the cloud.

- 3) *Fog and Cloud Are Mutually Beneficial*: Some functions are naturally more advantageous to be carried out in the fog while others in the cloud. Determining which functions should be carried out in the fog and how the fog should interact with the cloud will be key aspects of fog research and development.

Traditionally, services and applications are provided with large, centralized, expensive, and hard-to-innovate “boxes” such as the service gateways and packet data network gateways in the LTE core, large servers in a data center, and the core gateways and routers in a wide-area-network backbone. The traditional view is that the edge uses the core networks and data centers. The fog view is that the edge is part of the core network and a data center.

Table I outlines the main characteristics of fog and how it complements cloud.

A. Fog Architectural Characteristics and Advantages

A common denominator underlying fog is that fog distributes the resources and services of computation, communication, control, and storage closer to the users. A fog architectures may be fully distributed, mostly centralized, or somewhere in between. The fog architecture and the applications it supports may be virtualized but may also be implemented in dedicated hardware and software.

A fog architecture will allow the same application to run anywhere, reducing the need for specialized applications dedicated just for the cloud, just for the endpoints, or just for the edge devices. It will enable applications from different suppliers to run on the same physical platform without mutual interference. It will provide a common lifecycle management framework for all applications, offering capabilities for composing, configuring, dispatching, activating and deactivating, adding and removing, and updating applications. It will further provide a secure execution environment for fog services and applications. Fog will integrate with cloud to enable seamless end-to-end services.

Fog’s main advantages can be exemplified as CEAL.

- 1) *Cognition*: Awareness of client-centric objectives. A fog architecture, aware of customer requirements, can best determine where to carry out the computing, storage, and control functions along the cloud-to-thing continuum. Fog applications, being close to the end users, can be built to be better aware of and closely reflect customer requirements.
- 2) *Efficiency*: Pooling resources along the cloud-to-thing continuum. Fog can distribute computing, storage, and control functions anywhere between the cloud and the endpoint to take full advantage of the resources available along this continuum. It can also allow applications to leverage the otherwise idling computing, storage, and networking resources abundantly available on network edge and end-user devices such as tablets, laptops, smart home appliances, connected vehicles and trains, and network edge routers. Fog’s closer proximity to the endpoints will enable it to be more closely integrated with

TABLE I
MAIN CHARACTERISTICS OF FOG AND HOW IT COMPLEMENTS CLOUD

	Cloud	Fog
Location and Model of Computing	Centralized in a small number of big data centers.	Often distributed in many locations, potentially over large geographical areas, closer to users along the Cloud-to-Thing continuum. Distributed Fog nodes and systems can be controlled in centralized or distributed manners.
Size	Cloud data centers are very large in size, each typically contain tens of thousands of servers.	A Fog in each location can be small (e.g., one single fog node in a manufacturing plant or onboard a vehicle) or as large as required to meet customer demands. A large number of small Fog nodes may be used to form a large Fog system.
Deployment	Require sophisticated deployment planning.	While some Fog deployments will require careful deployment planning, Fog will enable ad-hoc deployment with no or minimal planning.
Operation	Operate in facilities and environments selected and fully controlled by Cloud operators. Operated and maintained by technical expert teams. Operated by large companies.	May operate in environments that are primarily determined by customers or their requirements. A Fog system may not be controlled or managed by anyone and may not be operated by technical experts. Fog operation may require no or little human intervention. May be operated by large and small companies, depending on size.
Applications	Support predominately, if not only, cyber-domain applications. Typically support applications that can tolerate round-trip delays in the order of a few seconds or longer.	Can support both cyber-domain and cyber-physical systems and applications. Can support significantly more time-critical applications that require latencies below tens of milliseconds or even lower.
Internet Connectivity and Bandwidth Requirements	Require clients to have network connectivity to the Cloud for the entire duration of services. Long-haul network bandwidth requirements grow with the total amount of data generated by all clients.	Can operate autonomously to provide uninterrupted services even no or intermittent Internet connectivity. Long-haul network bandwidth requirements grow with total the amount of data that need to be sent to the Cloud after being filtered by the Fog.

the end-user systems to enhance overall system efficiency and performance. This is especially important for performance-critical cyber-physical systems.

- 3) *Agility*: Rapid innovation and affordable scaling. It is usually much faster and cheaper to experiment with client and edge devices. Rather than waiting for vendors of large network and cloud boxes to initiate or adopt an innovation. Fog will make it easier to create an open market place for individuals and small teams to use open application programming interfaces, open software development kits (SDKs), and the proliferation of mobile devices to innovate, develop, deploy, and operate new services.
- 4) *Latency*: Real-time processing and cyber-physical system control. Fog enables data analytics at the network edge and can support time-sensitive functions for local cyber-physical systems. This is essential for not only stable control systems but also for the tactile Internet vision to enable embedded AI applications with millisecond reaction times as elaborated in the next subsection.

These advantages in turn enable new services and business models, and may help broaden revenues, reduce cost, or accelerate product rollouts as elaborated in the next subsection.

B. Fog Helps Address IoT Challenges

In particular, fog can provide effective ways to overcome many limitations of the existing computing architectures that rely only on computing in the cloud and on end-user devices. Table II shows, as an example, how fog can help address the IoT challenges we have discussed in Section II.

Proof-of-concept (POC) trials are demonstrating the business value and technology necessity of fog. For example, in late 2015, Cisco conducted a successful POC in Barcelona, where fog made smart city applications more cost-effective and manageable. Barcelona envisions deploying thousands of roadside cabinets throughout the city to optimize traffic management, energy management, and water and waste management. Before they could turn this vision into reality, the city faced two major challenges. First, the traditional way of adding new applications by adding dedicated new gateways and servers in every roadside cabinet is no longer feasible due to limited cabinet space. Second, the siloed applications have been using siloed application management systems, which made the system excessively expensive to deploy, operate, and maintain. Fog provided a solution. A single fog node provided a common platform at each cabinet for all services, and allowed applications from different suppliers to coexist without interfering with each other. It provided a unified platform to support networking, security, and lifecycle management for all applications, which reduced the systems costs and allowed application providers to focus on developing applications rather than providing specialized hardware and software to host and manage their applications.

C. Fog Enables New and Disruptive Business Models

Fog will enable new, and potentially highly disruptive, business models for computing and networking.

- 1) With fog, routers, switches, application servers, and storage servers will converge into fog nodes, with each fog node providing a common hardware and software

TABLE II
FOG PROVIDES EFFECTIVE WAYS TO ADDRESS IoT CHALLENGES

IoT Challenges	How Fog Can Help
Latency Constraints	Fog, performing data analytics, control, and other time-sensitive tasks close to end users, is the ideal and often the only option to meet the stringent timing requirements of many IoT systems.
Network Bandwidth Constraints	Fog enables hierarchical data processing along the Cloud-to-Things continuum, allowing processing to be performed where it can balance between application requirements and available networking and computing resources. This also reduces the amount of data that needs to be sent to the Cloud.
Resource-Constrained Devices	Fog can carry out resource-intensive tasks on behalf of resource-constrained devices when such tasks cannot be moved to the Cloud due to any reason, hence reducing these devices' complexity, lifecycle costs, and energy consumption.
Uninterrupted Services with Intermittent Connectivity to the Cloud	A local Fog system can operate autonomously to ensure non-interrupted services even when it has intermittent network connectivity to the Cloud.
New IoT Security Challenges	A Fog system can, for example, 1) act as the proxies for resource-constrained devices to help manage and update the security credentials and software on these devices, 2) perform a wide range of security functions, such as malware scanning, for the resource-constrained devices to compensate the limited security functionality on these devices, 3) monitor the security status of nearby devices, and 4) take advantage of local information and context to detect threats on a timely manner.

platform to support computing, networking, and storage. Such a transformation can significantly reshape the networking, server, and software industry landscape.

- 2) Fog-as-a-Service (FaaS) will enable new business models to deliver services to customers. Unlike the clouds that are mostly operated by large companies who can afford to build and operate huge data centers, FaaS will enable companies, big and small, to deploy and operate private or public computing, storage, and control services at different scales to meet the needs of a wide variety of customers.
- 3) Fog also provides a new way for network service providers to add value to customers in a new net-neutrality world. Consider, for example, the impact of the United States Federal Communications Commission (FCC) Title II ruling. The FCC vote in February 2015 to classify Internet services, including mobile services, as a "utility" under Title II regulatory

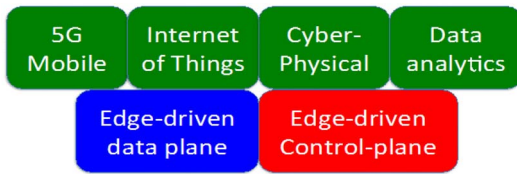


Fig. 1. Data plane and control plane of fog enable different applications.

mandate, may further push network innovation to the edge in the U.S. A new regulatory environment does not mean networks cannot be engineered and managed anymore, but we may need different vantage points of control: not from inside the network but from around the end users. For example, today network operators can pick which lane (WiFi, Macro-cellular, and Femtocell) a user device should be in. Since different lanes have different speeds and different payment system/amount, such practice may not be allowed any more in the U.S. Instead, we need new systems where each user device must choose which lane to be in for itself. As long as the government does not prohibit end-user choices, then we can run fog-based networking from the edge, through client/home-driven control/configuration.

IV. FOG USE CASE STUDIES

Architectural research and development asks the question of “who does what, at what timescale, and how to put the modules back together?” As an architecture, fog supports a variety of applications, including those typically associated with IoT and those often viewed as part of fifth-generation (5G) or data analytics and data management. Fog is an architecture for computing, storage, as well as for networking. In particular, fog network architecture consists of both data plane and control plane, each with a rapidly growing number of examples across protocol layers from the physical layer to the application layer.

1) Examples of Data Plane of Fog:

- a) pooling of clients idle computing/storage/bandwidth resources and local content;
- b) content caching at the edge and bandwidth management at home;
- c) client-driven distributed beam-forming;
- d) client-to-client direct communications (e.g., FlashLinQ, LTE direct, WiFi direct, and Air Drop);
- e) cloudlets and micro data-centers.

2) Examples of Control Plane of Fog:

- a) over the top (OTT) content management;
- b) fog-RAN: Fog driven RAN;
- c) client-based HetNets control;
- d) client-controlled cloud storage;
- e) session management and signaling load at the edge;
- f) crowd-sensing inference of network states;
- g) edge analytics and real-time stream-mining.

Data-plane of fog has been more extensively studied, e.g., [3]. In this section, we highlight a few particular cases that illustrate the potential and challenges of fog control plane,

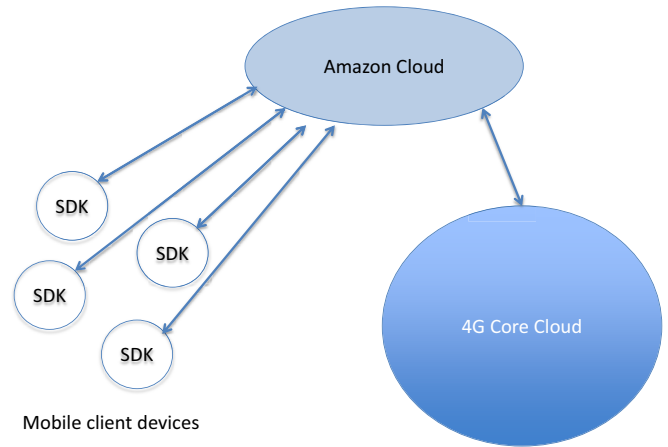


Fig. 2. SDK sitting inside clients can enable network inference and configuration. Crowd-sensing and byte-counting happen in fog, coordinated through a controller sitting in the Amazon cloud, bypassing the traditional reliance on network policy and configuration in the cellular core cloud.

such as the inference, control, configuration, and management of networks. We will also see that fog operates across a continuum spanning device, access, edge and more, and observe the collaboration between fog and cloud.

Case 1 [Crowd-Sensing LTE States (in Commercial Deployment)]: Through a combination of passive measurement (e.g., RSRQ), active probing (e.g., packet train), application throughput correlation and historical data mining, a collection of client devices may be able to, in real-time and useful accuracy, infer the states of an eNB such as the number of resource blocks used [4].

Case 2 [OTT Network Provisioning and Smart Data Pricing (in Commercial Deployment)]: [27] Fog directly leverages the “things” and phones instead, and removes the dependence on boxes-in-the-network altogether. With SDKs sitting behind apps on client devices, through tasks such as byte-counting, content tagging, location tracking, behavior monitoring, network services can be innovated much faster. In this case, the client SDKs collectively work through a controller (in the cloud as hosted say by Amazon) but bypass most of the cellular core network (a second type of cloud).

Case 3 [Client-Based HetNets Control (in 3GPP Standards)]: Coexistence of heterogeneous networks (e.g., LTE, femto, and WiFi) coexistence is a key feature in cellular networks today. Rather than through network operator control, each client can observe its local conditions and make decision on which network to join. Through randomization and hysteresis, such local actions may emerge globally to converge to a desirable configuration [5]. In the case of hybrid control of HetNets, the fog-cloud interface allows real-time network configuration be carried out by the clients themselves, while over longer timescale parameters like RAT stability attribute or hysteresis values can pass from the cloud (wireless core network) to the clients.

Case 4 [“Shred and Spread” Client-Controlled Cloud Storage (in Beta Trial)]: By decoupling massive cheap storage (in the cloud) from client side control of privacy (in the fog), we can achieve the best of both worlds. For example, by

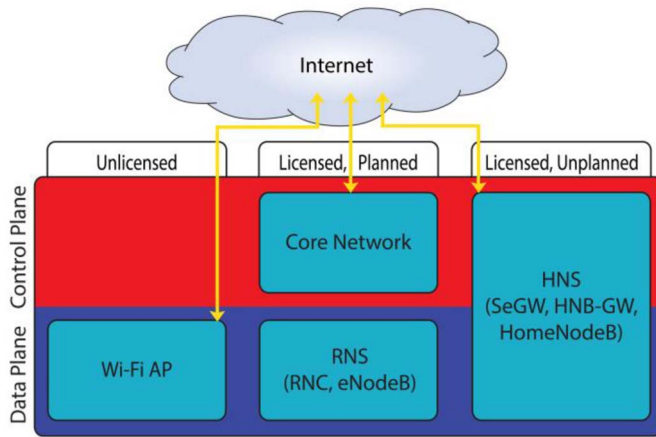


Fig. 3. Co-existence of heterogeneous networks may be managed in part by clients. Real time HetNets selection happens in fog, possibly aided by long-timescale parameter update from cloud [5].

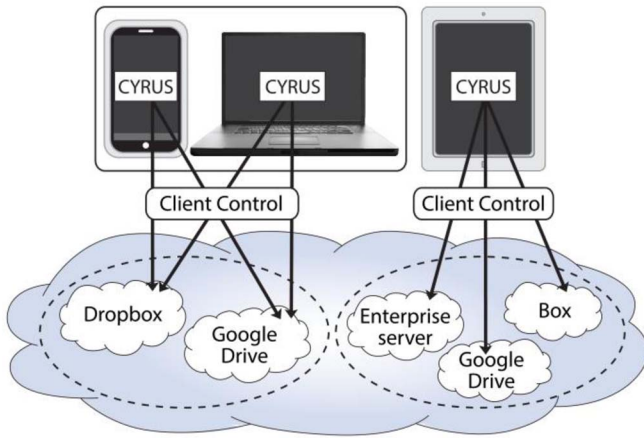


Fig. 4. Shred and spread (CYRUS project) shreds files in fog and spreads them over multiple clouds. Meta-data privacy control takes place in fog, while scalable storage remains in cloud [6].

shredding a file on the fog side and then spreading the bytes across multiple public clouds, in a client shim layer, of a given file across multiple cloud storage providers, it can be assured that privacy of the data is maintained even if encryption key is leaked by any given cloud provider [6].

Case 5 [Real-Time Stream Mining for Embedded AI (in Beta Trial)]: Consider virtual reality tasks associated with Google glass. Some of the information retrieval and computation tasks may be carried out on the glass (a “wearable thing”), some on the associated phone (a client device), some on the home storage (an edge device), and the rest in the cloud. An architecture of successive refinement may leverage all of these devices at the same time, with an intelligent division of labor across them [7].

Case 6 [Borrowing Bandwidth From Neighbors in D4D (in Beta Trial)]: When multiple devices belonging to the same person, to relatives or to employees of the same company are next to each other, one can ask the others to share their LTE/WiFi bandwidth by downloading other parts of the same file and transmitting, via WiFi direct, client to client [8].

Case 7 [Bandwidth Management at Home Gateway (in Beta Trial)]: By adapting the home set-top box/gateway, the limited broadband capacity is allocated among competing users and application sessions, according to each session’s priority and individual preferences. A prototype on a commodity router demonstrates a scalable, economical and accurate control of capacity allocation on the edge [9].

Case 8 [Distributed Beam-Forming (in Laboratory Demonstration)]: Fog can also happen in the physical layer, for example, by exploiting multiuser MIMO to improve throughput and reliability when a client can communicate with multiple WiFi access points. For uplink, we can use multiuser beam-forming so that the client can send multiple data streams to multiple APs simultaneously. For downlink, we can use interference nulling so that the client can decode parallel packets from multiple APs. These can be done entirely on the client side [10].

V. OPEN QUESTIONS AND RESEARCH CHALLENGES

As is typical of any emergent area of research and development, many themes in fog are not completely new, and instead are evolved versions of accumulated transformations in the past decade or two as follows.

- 1) Compared to P2P networks in the mid-2000s, fog is not just about content sharing (or data plane as a whole), but also network measurement, network management, service enablement, and real-time control of cyber-physical systems.
- 2) Compared to mobile *ad-hoc* network (MANET) research, fog will build upon much more powerful, diverse, and often off-the-shelf edge devices, applications, and end-to-end hierarchical networks enabled by broadband wireless and wired networks.
- 3) Compared to the generic edge-networking work in the past, fog adds a new layer of meaning to the end-to-end principle: in addition to optimizing among themselves, edge devices, collectively measuring and controlling the rest of the network, will collaborate with the cloud to enable end-to-end services along the cloud-to-things continuum.

Along with several other network architecture themes with longer histories, such as information-centric networks (ICNs), content-centric networks (CCNs), software-defined networks (SDNs), and network function virtualization (NFV), fog is revisiting the foundation of how to architect computing and networking: who does what and how to glue them back together.

- 1) *ICN/CCN*: Redefine functions (to operate on digital objects rather than just bytes).
- 2) *SDN/NFV*: Virtualize functions (through centralized control plane).
- 3) *Fog*: Relocate functions (closer to end users along the cloud-to-things continuum).

While fog does not have to rely on virtualization or to be information-centric, one can envision an information-centric and virtualized fog since these branches are complementary to each other and can be enablers for fog.

Fog also includes both mobile and wireline networks, and traverses edge, access and the wearables. As an important special case, supporting mobile edge computing inside an RAN will require many of the same functions of an end-to-end fog architecture to, for example, distribute, orchestrate, manage and secure the applications, and application enablement platforms. Fog, however, is broader than just supporting mobile edge computing. Fog is an architecture for distributing computing, storage, control, and networking services anywhere along the cloud-to-thing continuum, over and inside wireless and wireline networks, and supporting both mobile and wireline network applications.

As in any emergent area in its infant age, there is no shortage of challenging questions in fog, some of which continue from earlier study of P2P, MANET, and cloud, while others are driven by a confluence of recent developments in network engineering, user devices, and user experience.

The most fundamental question for fog in the coming decade is on where, when and how to distribute computation, communication, control and storage along the continuum from cloud to things. For example, how to decompose and re-compose computational tasks over a set of heterogeneously capable and variably available fog nodes wireless connected with bandwidth and energy constraints. Next, we discuss several categories of such fog research challenges.

A. Fog Interfaces With Cloud, Other Fogs, Things, and End Users

The fundamental question of architecture is “who does what, at what timescale, and how to put them back together?” In the case of fog, the question becomes:

- which tasks should go to the fog (e.g., those requiring real-time processing, end user objectives or low-cost leverage of idle resources);
- which go to the cloud (e.g., massive storage, heavy-duty computation, or wide-area connectivity);
- which go to the things;
- how the fog, the cloud, and the things should interact with each other.

The fog architectures should allow computing, storage, and networking tasks to be dynamically relocated among the fog, the cloud, and the things.

Therefore, the interfaces for fog to interact with the cloud, other fogs, and the things and users, as illustrated in Fig. 5, must: 1) facilitate flexible, and in some cases dynamic, relocation of the computing, storage, and control functions among these different entities; 2) enable convenient user access to fog services; and 3) allow efficient and effective lifecycle management of the system and services.

- 1) *Fog-to-Cloud Interfaces*: The fog-to-cloud interfaces will be needed to support fog-cloud collaborations to provide end-to-end services. It will support functions to, for example, allow:
 - a) fog to be managed from the cloud;
 - b) fog and cloud to send data to each other;
 - c) cloud to distribute services onto fog;
 - d) cloud to provide services to fog;

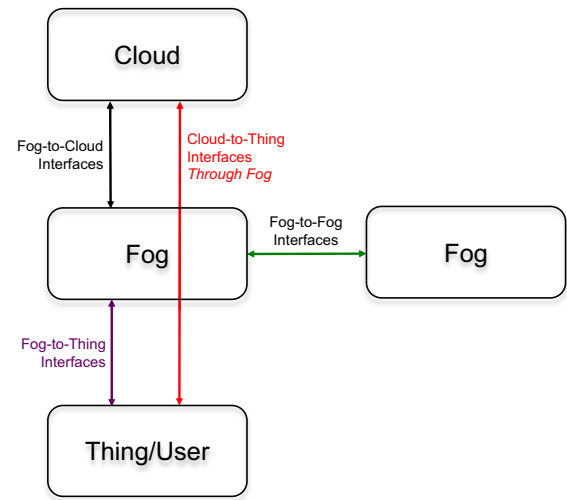


Fig. 5. Fog interfaces.

- e) cloud services to be provided through fog to things and end users;
- f) fog to provide services to cloud;
- g) fog and cloud to collaborate with each other to deliver end-to-end services.

It is essential to determine what information should be passed across the fog–cloud interface, the frequency and granularity of such information, and how the fog and the cloud should react to the information.

- 2) *Fog-to-Fog Interfaces*: Different fog nodes or systems may collaborate with each other to jointly support an application. For example, multiple fog systems can share the data storage and computing tasks for one or multiple users or applications. Different fog nodes or systems can also collaborate to serve as backups for each other. An important question is therefore how to design the interface and protocols to enable fog nodes in the same or different fog systems to collaborate.
- 3) *Fog-to-Thing/User Interfaces*: Fog will provide services to a wide range of end users and devices with widely varying capabilities. The fog-to-thing interface and fog-to-user interface will be essential to allow things and end users to access fog services in user-friendly, resource-efficient, and secure ways.

B. Fog-Enabled Edge and Access Networking

Fog can be used to support networking at the edge. For example, fog can provide services to help network edge devices and end-user devices (e.g., vehicles, drones, industrial and consumer robots, smartphones, and virtual reality gaggles) form local networks, providing temporary security credentials to these local devices to help them establish trustworthy communications, and act as local application servers and data storage servers for the edge networks. Some fog functions for supporting such edge networking may be implemented on the end-user devices. In such cases, how fog functions interface with the operating systems and hardware of the end-user devices becomes essential. More than just

using D4D for pooling idle edge resources as discussed in the previous sections, new protocol stacks for end-user devices to support fog-enabled edge networking may be needed.

C. Security

Fog presents new security challenges. Distributed systems are in general more vulnerable to attacks than centralized systems. While cloud operates in heavily protected facilities selected and controlled by cloud operators, fog often needs to operate in more vulnerable environments—where they can best meet customer requirements and often wherever users want them to be. Many fog systems will be significantly smaller than clouds (e.g., a fog node on a vehicle, in a manufacturing plant, or on a oil rig), and hence, may not have as much resources as the clouds to protect themselves. Furthermore, each fog system may not have the global intelligence necessary for detecting threats.

At the same time, however, fog's proximity to end users and locality on the edge enable it to help address certain new IoT security challenges as discussed in the previous sections. Fog can, for example, act as the first nodes for access control and traffic encryption, provide contextual integrity and isolation, serve as the aggregation and control points for privacy-sensitive data before the data leaves the edge, and act as the proxies of resource-constrained devices to carry out selected security functions for these resource-constrained devices.

D. Incentivization of Device Participation

In some IoT use cases, it is not too many un-trustworthy clients that create concerns but too few clients willing to participate. This can be the case when, for example, clients are expected to voluntarily contribute their computing or storage resources or to collaborate with each other to support applications. Market systems and incentive mechanisms will become useful.

E. Convergence and Consistency

Local interactions could lead to divergence, oscillation, and inconsistency of global system states, which are typical issues in distributed systems and can become more acute in a massive, under-organized, possibly mobile fog system with diverse capabilities, and potentially virtualized pool of resources shared unpredictably. Use cases in edge analytics and stream mining provide additional challenges on this recurrent challenge in distributed systems.

F. End-to-End Architectural Tradeoffs

Fog will create new opportunities for us to design end-to-end systems to achieve better tradeoffs between distributed and centralized architectures, between what stays local and what goes global, and between careful deployment planning and resilience through redundancy. Logical fog system topologies, statically or dynamically established, over the same underlying physical fog network can be used to support a spectrum of architectures from completely centralized to fully distributed.

To address the above challenges, we need both of the following:

- 1) fundamental research, across networking, device hardware/OS, pricing, human-computer interface, and data science;
- 2) industry-academia interactions, as exemplified in the Open Fog Consortium (OpenFog), a global, nonprofit consortium launched in November 2015 with founding members from ARM, Cisco, Dell, Intel, Microsoft and Princeton University Edge Laboratory.

VI. CONCLUDING REMARKS

Fog is starting to reshape the future landscape of multiple industries, driving innovation across the entire industry food chain, including the following.

- 1) End user experience providers (e.g., GE, Toyota, Sony, Walmart, etc.).
- 2) Network operators (e.g., AT&T, Verizon, Comcast, etc.).
- 3) Network equipment vendors (e.g., Cisco, Nokia, Ericsson, Huawei, etc.).
- 4) Cloud service providers (e.g., VMWare, Amazon, etc.).
- 5) System integrators (e.g., IBM, HP, etc.).
- 6) Edge device manufacturers (e.g., Linksys, Samsung, etc.).
- 7) Client and IoT device manufacturers (e.g., Dell, Microsoft, Apple, Google, etc.).
- 8) Computer chip suppliers (e.g., Intel, ARM, Qualcomm, Broadcom, etc.).

The past 15 years have seen the pendulum swinging toward “click.” Now it has started to swing back closer to the “brick,” pointing to a co-existence of fog and cloud. Cloud has the advantages for massive storage, heavy duty computation, global coordination and wide-area connectivity, while fog will be useful for real time processing, rapid innovation, user-centric service and edge resource pooling. 2016 is an interesting year to start systematically exploring what fog might look like and the differences it will bring to the world of networking and computing in the next 15 years.

ACKNOWLEDGMENT

The authors are grateful for the inspiring conversations with colleagues at Princeton Edge Lab and Cisco, as well as with many colleagues in industry and academia, especially, Flavio Bonomi, Russell Hsing, Bharath Balasubramanian, Sangtae Ha, Junshan Zhang, Raj Savor, John Smee, Chonggang Wang, and representatives of many member companies and universities in the Open Fog Consortium.

REFERENCES

- [1] M. Chiang, S. H. Low, A. R. Calderbank, and J. C. Doyle, “Layering as optimization decomposition: A mathematical theory of network architectures,” *Proc. IEEE*, vol. 95, no. 1, pp. 255–312, Jan. 2007.
- [2] F. Bonomi, R. Milito, P. Natarajan, and J. Zhu, “Fog computing: A platform for Internet of Things and analytics,” in *Big Data and Internet of Things: A Roadmap for Smart Environments*. Cham, Switzerland: Springer, Mar. 2014, pp. 169–186.
- [3] M. Satyanarayanan, P. Bahl, R. Caceres, and N. Davies, “The case for VM-based cloudlets in mobile computing,” *IEEE Pervasive Comput.*, vol. 8, no. 4, pp. 14–23, Oct./Dec. 2009.

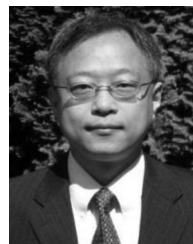
- [4] A. Chakraborty, V. Navda, V. N. Padmanabhan, and R. Ramjee, "Coordinating cellular background transfers using loadsense," in *Proc. Mobicom*, Miami, FL, USA, 2013, pp. 63–74.
- [5] E. Aryafar, A. Keshavarz-Haddad, M. Wang, and M. Chiang, "RAT selection games in HetNets," in *Proc. IEEE INFOCOM*, Turin, Italy, 2013, pp. 998–1006.
- [6] J. Y. Chong, C. Joe-Wong, S. Ha, and M. Chiang, "CYRUS: Towards client-defined cloud storage," in *Proc. EuroSys*, Bordeaux, France, 2015, Art. no. 17.
- [7] L. Canzian and M. van der Schaar, "Realtime stream mining: Online knowledge extraction using classifier networks," *IEEE Netw.*, vol. 29, no. 5, pp. 10–16, Sep./Oct. 2015.
- [8] X. Chen, B. Proulx, X. Gong, and J. Zhang, "Social trust and social reciprocity based cooperative D2D communications," in *Proc. ACM Int. Symp. Mobile Ad Hoc Netw. Comput. (MOBIHOC)*, Bengaluru, India, Jul. 29–Aug. 1, 2013, pp. 187–196.
- [9] F. M. F. Wong, S. Ha, C. Joe-Wong, Z. Liu, and M. Chiang, "Mind your own bandwidth: Adaptive traffic management on network edge," in *Proc. IEEE IWQoS*, Portland, OR, USA, 2015.
- [10] Y. Du, E. Aryafar, J. Camp, and M. Chiang, "iBeam: Intelligent client-side multi-user beamforming in wireless networks," in *Proc. IEEE INFOCOM*, Toronto, ON, CA, 2014, pp. 817–825.
- [11] M. Weiner, M. Jorgovanovic, A. Sahai, and B. Nikolić, "Design of a low-latency, high-reliability wireless communication system for control applications," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Sydney, NSW, Australia, 2014, pp. 3829–3835.
- [12] R. Kelly, *Internet of Things Data to Top 1.6 Zettabytes by 2022*. Accessed on Apr. 7, 2016. [Online]. Available: <https://campustechnology.com/articles/2015/04/15/internet-of-things-data-to-top-1-6-zettabytes-by-2020.aspx>
- [13] L. Mearian, *Self-Driving Cars Could Create 1GB of Data a Second*. Accessed on Apr. 7, 2016. [Online]. Available: <http://www.computerworld.com/article/2484219/emerging-technology/self-driving-cars-could-create-1gb-of-data-a-second.html>
- [14] N. Cochrane, (Mar. 23, 2010). *US Smart Grid to Generate 1000 Petabytes of Data a Year*. Accessed on Apr. 7, 2016. [Online]. Available: <http://www.itnews.com.au/news/us-smart-grid-to-generate-1000-petabytes-of-data-a-year-170290#ixzz458VaITi6>
- [15] K. Stouffer, J. Falco, and K. Scarfone, "Guide to industrial control systems (ICS) security," Nat. Inst. Standards Technol. (NIST), U.S. Dept. Commerce, Washington, DC, USA, Special Pub. 800-82, Jun. 2011.
- [16] W. Ashford, (Oct. 15, 2014). *Industrial Control Systems: What are the Security Challenges?* Accessed on Jan. 28, 2016. [Online]. Available: <http://www.computerweekly.com/news/2240232680/Industrial-control-systems-What-are-the-security-challenges>
- [17] Bureau of Transportation Statistics, U.S. Dept. Transp., Washington, DC, USA. Accessed on Mar. 2, 2016. [Online]. Available: http://www.rita.dot.gov/bts/sites/rita.dot.gov/bts/files/publications/national_transportation_statistics/html/table_01_26.html_mfd
- [18] G. Gan, Z. Lu, and J. Jiang, "Internet of Things security analysis," in *Proc. Int. Conf. Internet Technol. Appl. (iTAP)*, Wuhan, China, Aug. 2011, pp. 1–4.
- [19] N. Falliere, L. O. Murchu, and E. Chien, "W32.stuxnet Dossier," Symantec Security Response, Ver. 1.4, Mountain View, CA, USA, Feb. 2011.
- [20] K. Zetter, *Countdown to Zero Day: Stuxnet and the Launch of the World's First Digital Weapon*. New York, NY, USA: Crown, 2014.
- [21] *Stuxnet*. Accessed on Mar. 2, 2016. [Online]. Available: <https://en.wikipedia.org/wiki/Stuxnet>
- [22] L. Delgrossi and T. Zhang, *Vehicle Safety Communications: Protocols, Security, and Privacy*. Hoboken, NJ, USA: Wiley, 2012.
- [23] T. Zhang, H. Antunes, and S. Aggarwal, "Defending connected vehicles against malware: Challenges and a solution framework," *IEEE Internet Things J.*, vol. 1, no. 1, pp. 10–21, Feb. 2014.
- [24] T. Zhang, H. Antunes, and S. Aggarwal, "Securing connected vehicles end to end," in *Proc. SAE World Congr. Exhibit.*, Detroit, MI, USA, Apr. 2014.
- [25] R. Chen, L. Wei, H. Zou, and M. Zhai, "A TCM-based remote anonymous attestation protocol for power information system," in *Proc. Int. Power Electron. Mater. Eng. Conf. (IPEMEC)*, Dalian, China, May 2015.
- [26] A. Francillon, Q. Nguyen, K. B. Rasmussen, and G. Tsudik, "A minimalist approach to remote attestation," in *Proc. Conf. Design Autom. Test Europe (DATE)*, Dresden, Germany, 2014, pp. 1–6.
- [27] S. Ha, S. Sen, C. Joe-Wong, Y. Im, and M. Chiang, "TUBE: Time-dependent pricing for mobile data," in *Proc. ACM SIGCOMM*, Helsinki, Finland, Aug. 2012, pp. 247–258.



Mung Chiang (S'00–M'03–SM'08–F'12) is the Arthur LeGrand Doty Professor of Electrical Engineering with Princeton University, Princeton, NJ, USA. His book *Networks: Friends, Money and Bytes* and online course reached 250 000 students since 2012. He founded the Princeton Edge Laboratory in 2009, which bridges the theory-practice gap in edge networking research by spanning from proofs to prototypes. He co-founded a few startups in mobile pricing, IoT, big data areas, and the Open Fog Consortium. He is the Director

of Keller Center for Innovations in Engineering Education in Princeton University and the inaugural Chairman of the Princeton Entrepreneurship Council.

Prof. Chiang was the recipient of the 2013 Alan T. Waterman Award, the highest honor to U.S. young scientists and engineers.



Tao Zhang (F'00) received the B.S. and M.S. degrees in electrical engineering from Northern Jiaotong University, Beijing, China, and the Ph.D. degree in electrical and computer engineering from the University of Massachusetts, Amherst, MA, USA.

He joined Cisco Systems, San Jose, CA, USA, as the Chief Scientist for Smart Connected Vehicles, in 2012. Since then, he has also been leading initiatives to develop strategies, architectures, technology, and eco-systems for the Internet of Things and fog computing. From 1995 to 2012, he was with Telcordia Technologies (formerly Bell Communications Research or Bellcore), Piscataway, NJ, USA, where he was a Chief Scientist and the Director of Mobile and Vehicular Networking. For over 25 years, he has been in various technical and executive positions, directing research and product development. He holds over 50 U.S. patents and co-authored two books *Vehicle Safety Communications: Protocols, Security, and Privacy* (Wiley, 2012) and *IP-Based Next Generation Wireless Networks* (Wiley, 2004).

Dr. Zhang is a Co-Founder and a Board Director of the Open Fog Consortium and the CIO and a Board Governor of the IEEE Communications Society. He was the Founding Board Director of the Connected Vehicle Trade Association.