

Toward a Global Data Infrastructure

Nitesh Mor, Ben Zhang, John Kolb, Douglas S. Chan, Nikhil Goyal, Nicholas Sun,
Ken Lutz, Eric Allman, John Wawrzynek, Edward Lee, John Kubiawicz
University of California, Berkeley

Abstract

The Internet of Things (IoT) represents a new class of applications that can benefit from cloud infrastructure. However, directly connecting smart devices to the cloud has a number of disadvantages and is unlikely to keep up with either the growing speed of the IoT or the diverse needs of IoT applications.

We explore these disadvantages and argue that fundamental properties of the IoT prevent the current approach from scaling. What is missing is a well-architected system extending functionality of the cloud and providing seamless interplay among the heterogeneous components closer to the edge in the IoT space. We argue that raising the level of abstraction to a data-centric design—focused around the distribution, preservation and protection of information—better matches the IoT. We present early work on such a distributed platform, called the Global Data Plane (GDP), and discuss how it addresses the problems with the cloud-centric architecture.

1 Introduction

The market has seen an explosion in the number of smart devices. These latest devices offer rich interactivity by connecting to computing platforms and services. Fueled by the growth of *Internet* connectivity and the augmentation of everyday *things*, this shift is commonly referred to as the Internet of Things (IoT). With the ever-growing proliferation of cloud computing and storage services, even novice users can start deploying sensors and actuators with minimal invest-

ment in infrastructure. However, issues of privacy, security, scalability, latency, bandwidth availability, *etc.* that already are a challenge for web-applications, are exacerbated in the IoT space because of the fundamental differences between IoT and web services (see Sec. 2).

In this article, we analyze the shortcomings of existing architecture by explaining the fundamental differences between IoT applications and web services. Our analysis suggests a need for a new layer of abstraction for the IoT—one that more naturally fits the requirements of IoT applications while exploiting the underlying computing platforms that enable the IoT (like the cloud, the Fog³ and gateways). Although influenced by the needs of IoT, there is no reason this layer of abstraction can not be used for other scenarios.

Our new abstraction is centered around data. It is focused on the transport, replication, preservation, and integrity of streams of data while enabling transparent optimization for locality and quality of service. We call the resulting infrastructure the Global Data Plane (GDP). Its foundation is the concept of a single-writer append-only log coupled with location-independent routing, overlay multicast and higher level interfaces such as common access APIs (see Sec. 3). We also discuss how a communication and storage platform can enable better security by reducing attack surface of potentially vulnerable end-devices. Since this is an ongoing effort, we focus mainly on the design experience with GDP thus far.

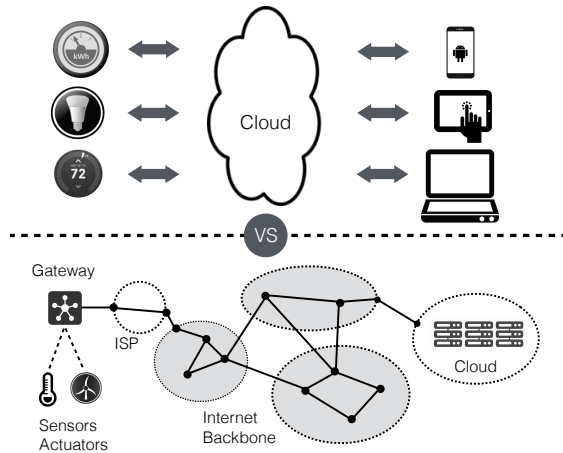


Figure 1: Although applications usually view the cloud as the center of all connected devices (*upper diagram*), in reality the cloud is usually on the edge of the Internet backbone, just like other devices (*lower diagram*).

2 Pitfalls with Today’s Approach to IoT

Looking at recent trends, IoT applications fall into two general categories:

I. Ambient data collection and analytics: These applications involve sensors installed in buildings, in cities, and on humans themselves. Normally, data is not immediately inspected and the collected data is later processed for analytics.

II. Real-time applications with low-latency requirements: These applications could be either autonomous systems with tight control loops (*e.g.* robots taking actions based on sensors), or reactive environments with humans in the loop. Unpredictable latencies of cloud-based solutions make this challenging.

It is not unusual to have same infrastructure be used for both categories of IoT applications. Performing actuation in real time using local sensor data, but with global knowledge of long term trends to optimize on various parameters is a common practice. With this context, relying entirely on the cloud for IoT applications has the following concerns:

1. The Cloud is Suboptimal : Application developers view the cloud as an interconnection hub for smart devices. However, from a networking point of view, the cloud is on the edge of the network, leading to unpredictable latencies (see Fig. 1). For

web-applications, a centralized hub enables amortization of resources and enables economies of scale (*e.g.* by caching popular resources, reduced management overhead, *etc.*), and occasional latencies are often ignored in favor of such cost benefits. IoT applications, on the other hand, are highly specific in their execution pattern and do not necessarily benefit from such centralization.

2. Security and Privacy: Sensors implanted in our surrounding environment may collect extremely sensitive information. As a centralized resource out of users’ control, the cloud presents an ever-present opportunity to violate privacy, already a luxury, and threatened further by the IoT. With appropriate encryption techniques, data stored in cloud can be protected from unauthorized use. However, an end-to-end secure data flow is hard to achieve because most applications can’t process encrypted data. Further, there is information leak from side-channels based on access patterns. Even with the state of the art (homomorphic encryption, secure hardware containers, *etc.*), absolute end-to-end security is a challenging technical problem. Giving users some control over where applications execute is a more general solution, especially when low latency requirements also dictate so.

3. Quality of Service: Guarantees on latency and availability are hard to realize. Web users tolerate variable latency and occasional loss of web services. In contrast, the temporary unavailability of sensors or actuators within IoT applications will directly impact the physical world. While significant engineering effort has been put into improving the availability and latency profile of the cloud, such efforts are stymied by operator error, software bugs, DDoS attacks, and normal packet-to-packet variations from wide-area routing. Further, the Internet connection to peoples’ homes is far from perfect even in developed world; this situation is worse in developing countries.

4. Durability Management: Some sensor data is ephemeral, while other data should be durable against global disasters. For ephemeral data, there is no effective way of verifying the data has been completely destroyed because the cloud is out of the user’s control. Control over durability is closely related to control over data in general: making sure that users retain the control and ownership over their

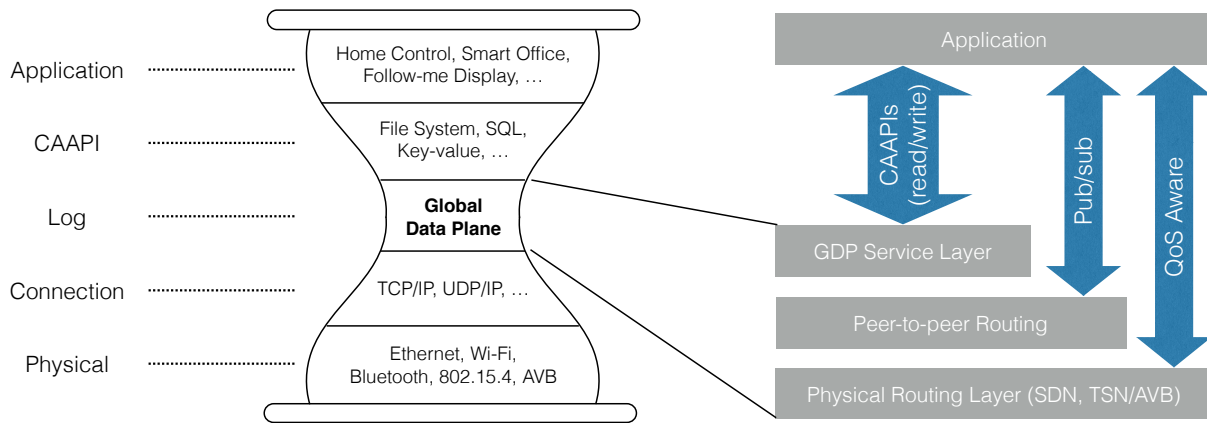


Figure 2: The Global Data Plane (GDP) operates above the network level and offers Common Access APIs (CAAPIs) to applications rather than raw packet routing. We argue that this abstraction is appropriate for IoT applications both in the cloud and in distributed infrastructure.

data rather than service providers.

3 GDP: A Data-Centric Proposal

In contrast to the existing cloud-centric model, we argue for the Global Data Plane (GDP), a data-centric abstraction focused around the distribution, preservation, and protection of information. It supports the same application model as the cloud, while better matching the needs and characteristics of the IoT by utilizing heterogeneous computing platforms such as small gateway devices, moderately powerful nodes in the environment and the cloud, in a distributed manner. The key mechanism for data storage and communication in GDP is the secure, single-writer log, which we describe in more detail later. As shown in Fig. 2, this log interface of the GDP provides a new “narrow waist” upon which applications are constructed.

3.1 System Overview

The concept of a log is central to the GDP. As the name suggests, a log is a time-series append-only data-structure addressed using a flat 256-bit identifier, called a *GDP-name*. A log is an *authenticated data structure*¹⁰ stored on potentially untrusted infrastructure. Logs are lightweight, durable, potentially distributed over multiple physical machines,

and don’t have a fixed location but rather are migrated as necessary to meet the locality, privacy, or QoS needs of applications. Logs are single-writer but support multiple simultaneous readers—either through random access (pull-based) or subscription (push-based). Not only logs, but other entities in the GDP (clients, log-servers, etc) have flat 256-bit GDP-names.

Clients in the GDP are entities that read from or write to logs—this includes sensors that generate data, actuators that consume data, gateway devices, and various software entities in-between that process data by reading it from an input log and writing it to an output log (see Fig. 3). We provide an event-driven communication library for interaction with logs. Moderately powerful gateway devices (smartphones, Raspberry Pi, *etc.*) are a particularly important case: they connect to the GDP on behalf of sensors/actuators limited in computation or communication capabilities and enable wide adoption with minimal changes in end-devices. We assume that clients have cryptographic capabilities, including key-storage (for encryption/decryption and signatures). Signatures are used for verifying the origin, authenticity and integrity of data flow and control commands. Encryption is used wherever necessary to provide data secrecy. We describe this in more detail in Sec. 3.4.1.

Logs are physically stored on *log-servers*—these log-servers can be small ubiquitous devices in homes,

local servers in an enterprise, or powerful cloud based servers backed by existing cloud storage systems. *GDP-routers* are the routing elements that provide *location-independent routing* in this large, 256-bit address space using an overlay network. We have a preliminary version of GDP-router implemented in Click⁶ on top of TCP/IP.

Our design is guided by the goal that a client should be able to operate without a single-point of trust in log-servers and GDP-routers. We hope to achieve this using a combination of cryptographic operations, trusted hardware and secure multi-party computation.

We also have a notion of a *Control Plane*—a set of services and applications that provide policy enforcement using mechanisms provided by the GDP. As an example, GDP makes sure that the logs are durable by ensuring replication across domains guided by a control plane replication service. The control plane and the GDP are closely integrated, yet have well specified boundaries.

3.2 The Log Interface

The majority of sensors and actuators in IoT can easily be represented by a stream of data, hence a queue-like interface for storing this streaming data seems to be the most obvious choice. However, the life-span of this data is application dependent. A log interface provides a wide range of options: logs could potentially be truncated for ephemeral data, or replicated widely for long lived data.

3.2.1 Properties of a Log

Logs are *append-only*; already existing data in a log is *read-only* and can be securely replicated and validated through cryptographic hashes and signatures. A *record* is the unit of read/write to a log; a log is essentially an ordered list of records. In addition, each log has immutable metadata created at the time of log-creation. For each log, our current design exposes *append*, *read* and *subscribe* APIs. Logs are single-writer, thus enabling serialization of records at client side. This implies that each sensor has its own log, however, aggregated logs representing more than one sensor can be created by reading multiple logs and writing back to another log. A single-writer

log is minimal but complete interface that could be used to build richer interfaces.

Even though a log is *append-only*, log-truncation policies can be specified on a per log basis. Log-truncation marks the data older than a specific threshold *safe to delete* by the infrastructure. Guaranteed destruction of data by a remote entity is practically impossible. However, since the GDP allows clients to specify where data is stored, it is much easier to control the longevity of data than in an exclusively cloud-based approach. Further, a client can encrypt a range of data with a unique key, and destroy the decryption key for data older than a certain threshold while maintaining data-integrity related invariants of the log data-structure.

A log is created by a client by issuing a signed *create-request*, which contains metadata including the public signature key of the designated writer. The create-request gets routed through a series of control plane services, the control plane checks whether the client is authorized to create a new log or not, allocates resources for this newly created log, sets up replication, etc.

Write access control is performed by the log-servers by validating signature on append operations against the designated writer's public signature key, while read access control is implemented by encrypting the payload and selectively sharing the decryption key. Since signatures remain with data, a malfunctioning or malicious log server is unable to fabricate data. More details on this are in Sec. 3.4.1.

In addition, a variety of basic control plane services could be used for making a log more functional. A replication service could set up multiple replicas of a log based on higher level policy decisions, such as level of durability, geographic span, log-truncation policies, etc. on a per log basis. A directory service could be used to associate human-readable names to 256-bit GDP-names on an organization level, or at a user level.

3.2.2 Benefits of a Log Interface

A log interface makes dumb sensors and actuators significantly more functional. Low-power sensors usually only generate data, but can not answer any queries. If data values are written to a log by such sensors, the log can be used as a proxy that supports

a much richer set of queries, especially for historical data. A subscription to such a log provides latest sensor values in almost real time, thus virtualizing the sensor in some sense.

Actuators, on the other hand, usually need to maintain some kind of access control—by physical isolation, some authentication method, or a combination of both. Instead, if an actuator were to subscribe to an actuation log to read the actuation commands, access control can be implemented at the log level. This makes actuator design simpler and avoids the pitfalls of ad-hoc authentication mechanisms hastily put together by hardware vendors.

Further, there is no need to expose the physical devices with potentially questionable standards of software security to the entire world, while still being able to connect things together. This is especially important since it takes the burden of implementing security off the device vendors’ shoulders. All a device manufacturer has to do is publish data to a log (in case of a sensor), or subscribe to a log (in case of an actuator). This greatly reduces the attack surface, because GDP as a platform implements security best practices.

Representing sensors and actuators with logs separates policy decisions from mechanisms, enabling cleaner application designs. Applications can be built on top of GDP by interconnecting globally addressable log streams, rather than by addressing devices or services via IP addresses. Further, with applications running inside containers (Docker, Unikernels, Intel SGX enclaves, *etc.*), forcing data-flows in and out of the container through logs enables any filtering at the log-level (*e.g.* access-control).

And last, but not the least, the “narrow waist” provided by globally addressable logs avoids stove-piped solutions and provides for a heterogeneous hardware infrastructure.

3.2.3 Beyond a Log Interface: Common Access APIs (CAAPIs)

Although a log abstraction shelters developers from low-level machine and communication primitives, many applications are likely to need more common APIs or data structures. In fact, logs are sufficient to implement any convenient, mutable data storage repository. Thus, Fig. 2 shows a Common Access

API (CAAPI) layer on top of the GDP. A CAAPI can present a key-value store, file system or database interface. Since logs serve as the ground truth, the benefit of consistency, durability, scalability and availability are carried over to CAAPIs for free.

3.3 Flat Address Space

As mentioned earlier, we use 256-bit long flat addresses for naming things in the GDP. This is true not only for logs, clients, log-servers, but also for control plane services and applications. In particular, logs are named with a 256-bit identifier which may be derived from a cryptographic hash of the owner’s public key and meta-data.

This large address space allows us to employ *location-independent routing* that can better match the goals of flexible placement, controllable replication and mobility to optimize for latency, QoS, privacy and durability. Following a variety of placement and replication policies, GDP places logs within the infrastructure and advertises the location of these logs to the underlying routing layer.

GDP-routers take the burden of performing and optimizing this location-independent routing through an overlay network that uses a combination of Distributed Hash Table (DHT) technology and selective routing. DHT addresses the challenges of scalability with the sacrifice of an increased number of overlay hops. Important routes can be optimized by pushing routing entries into an underlying routing layer, possibly using Software-Defined Networking (SDN) routers when available. Fig. 2 shows this layering.

Based on the interaction between GDP-routers and control-plane, logs could be migrated and routing topology altered dynamically. In addition, multicast trees can be built on top of the overlay network^{2,12} to efficiently serve multiple subscribers. Not only migration, but location-independent flat-addressing also enables replication; a log (or service) can have multiple read-only copies spread throughout the network. Single-writer append-only design makes a log mostly read-only (except for the active head), which makes for easy replication with simpler concurrency issues. These read-only replicas act very much like a Content Delivery Network (CDN) and provide for fault-tolerance. In case of network events such as flash crowds or targeted bandwidth saturation attacks,

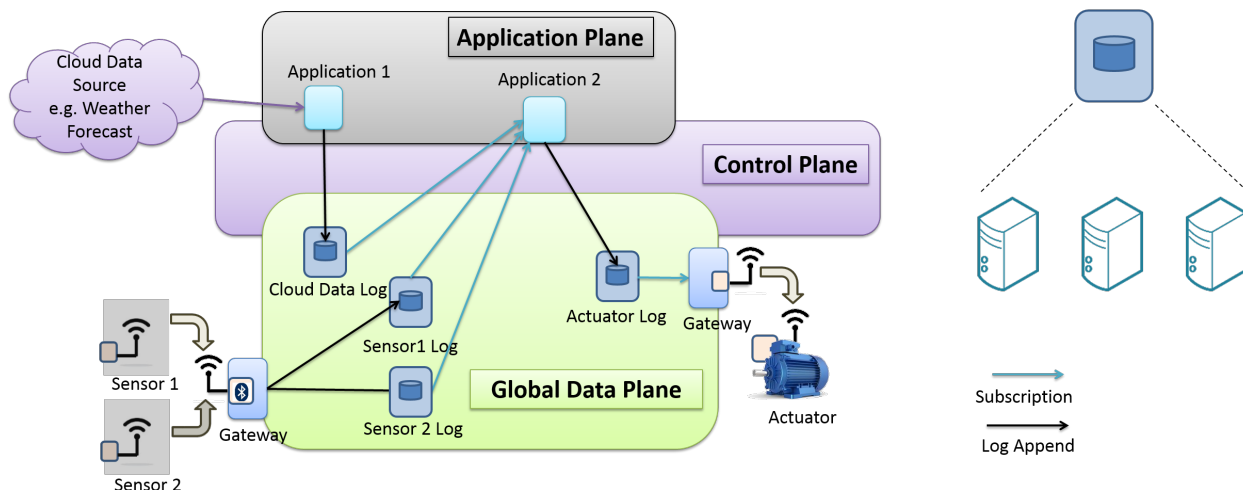


Figure 3: An IoT application uses GDP to combine heterogeneous data-streams from local environmental sensors and from a cloud source to actuate a device. Instead of direct communication with devices, the application uses the “narrow waist” (logs) provided by the GDP. Even though a log is represented as a single data-stream to an application, internally it can be distributed and replicated over multiple physical log servers to achieve locality and durability.

the traffic is distributed over multiple replicas rather than a single point of failure. In addition, since low-power devices (sensors/actuators) are exposed to the world as a virtual device as a log, any potential bandwidth saturation is limited to the GDP and does not turn into a battery exhaustion attack.

3.4 Challenges

In this section, we describe the major challenges that we faced in the design of GDP, and that are of concern to a general IoT framework.

3.4.1 Security and Privacy

As we have outlined in Sec. 2, data security and privacy is more important than ever, given the pervasive nature of devices and actuators. In the GDP, we design our security and privacy mechanisms and policies by focusing on the “narrow waist” provided by the logs.

Logs are stored on potentially untrusted log-servers. Hence it is important that we do not rely on a single log-server to provide data-integrity. An adversarial log-server could try to tamper with existing data in the logs it stores, or may not perform appropriate checks on access-control and accept writes from unauthorized writers, or maliciously re-order append operations received from a legitimate writer.

We address data-integrity and write-access control challenges using a combination of signatures and Merkle trees in our single-writer log model; a writer signs each append operation with a signature key and performs record-ordering on the client side by including a hash-pointer to the previous record in the signed content. The public signature key for the writer is included at the time of creation in the *create-request*, which itself is signed. All that a log-server has to do is to perform a signature-validation against this well-specified public key for any new append operation it receives. Any accidental or malicious behavior by a log-server results in invalid signatures or a broken chain of hash-pointers, and can be detected by a reader.

Globally addressable logs are a significant privacy concern if any unauthorized reader could read data at will. We envision encryption to be the mechanism for providing data secrecy. GDP does not assume any structure on the data being written to a log, enabling applications to encrypt data before handing it to GDP. This enforces the minimum trust philosophy by putting trust in cryptographic constructs rather than potentially buggy software running on untrusted servers. Read access-control is managed by the application by appropriate sharing of the decryption keys.

A secondary concern is exposing encrypted data

to adversaries who may analyze timing or data size at will. To address these concerns, we envision a collaboration between policy (at the control plane) and routing (within the GDP) to help mitigate this problem by controlling the placement of data logs and path of updates.

3.4.2 Key-management

Since security and privacy in the GDP relies upon encryption, key management will be of paramount importance. Although advanced schemes for key management are still under consideration, we support a basic key management scheme as follows.

Each user maintains an encrypted wallet and an unencrypted public-key registry, both backed by logs. The user-supplied root key is used to decrypt the wallet, which contains the secret keys necessary to sign requests and secret or symmetric keys necessary to decrypt log entries.

Granting read access amounts to sending a bit-string over a tamper-proof channel (a log) to a remote entity; this bit string is the necessary decryption key that is in turn encrypted using the public key of the remote entity. As a very simple example: Alice wants to share a log L with a set of users. Alice creates the log L including the name of an "access control log" A in L 's metadata. Alice then encrypts the contents of L with a symmetric key K and appends versions of K to A , each encrypted using the public key of the users who should have access.

This scheme works for simple and static data sharing scenarios. Slightly complicated but efficient hierarchical key management schemes can be created based on application requirements. In the extreme case, a compute service in a trusted environment could be designated as the only reader, with that service managing read access control lists in more traditional ways.

3.4.3 Routing Overhead

A flat address space offers several benefits, but a naïve overlay implementation can severely affect the performance, especially round-trip latencies. However, reasonably dense locality-aware distributing routing frameworks limit relative delay penalty (ratio between the distance traveled between two points

using an overlay and the minimum distance between the two points) of using a DHT to $2 - 3$.¹² We propose to use locality to overcome the overlay performance penalty; *e.g.* round-trip latency to a close-by/local node via overlay ($10 - 20ms$) is still better than round-trip latency to a cloud-datacenter ($> 50ms$). Further, as mentioned in Sec. 3.3, with support for optimizing flows from the underlying networks (SDN, AVB, *etc.*), long term communication to/from logs can achieve better performance.

4 Related Work

Other efforts exist to address the challenges of IoT. Ericsson's Capillary networks¹ and Cisco's Fog Computing³ provide computing resources closer to the edge of the network. We believe that our arguments strengthen the need for fog-like computing platforms and our GDP architecture can both leverage and simultaneously enable such platforms. Also relevant are systems such as EdgeComputing from Akamai,⁴ and Cloudlet.⁸ In these architectures, the role of servers is to be intelligent gateways or proxies for data flowing into and from the cloud. Support for an entirely decentralized data storage and delivery platform is apparently absent.

Our data-centric design hails from Oceanstore⁷ and shares a number of goals with Named Data Networking (NDN)¹¹ and MobilityFirst,⁹ but our focus on the IoT application space leads to a number of important design differences. Among other things, push based communication—such as from sensors to logs and from logs to consumers—represents a communication style utilized extensively in the IoT space and deemphasized in NDN. MobilityFirst shares the fundamental design principle of a flat-address (GUID), however the emphasis is primarily on communication. In contrast, GDP provides a higher level log abstraction supported by underlying communication primitives.

A few of our design decisions are similar to Bolt:⁵ single-writer time-series data, chunking for performance, efficient data sharing, policy-driven storage and data confidentiality/integrity. However, Bolt takes the cloud approach where the pitfalls in Sec. 2 are unavoidable.

5 Conclusions

A prototype version of the GDP has been deployed within our own environment and has been running on a few servers since early 2015, however it is still a work in progress. Our design for the GDP is not yet bullet-proof and our initial implementation has not withstood the test of wide-scale deployment. Nonetheless, we believe that the core concepts of GDP overcome the pitfalls mentioned in Sec. 2 in the following way: the single-writer, append-only log models sensor data more accurately; integrity and authentication by design provides better privacy and security; the distributed nature with peer-to-peer technology makes scalability possible; explicit separation of policy from mechanism enables better control on level of durability for end users; and finally, latency, bandwidth and QoS guarantees are enabled by the integration of the cloud and the local infrastructure.

6 Acknowledgments

This work was supported in part by the Ubiquitous Swarm Lab at UC Berkeley. This work was also supported in part by TerraSwarm, one of six centers of STARnet, a Semiconductor Research Corporation program sponsored by MARCO and DARPA.

References

- ¹ Capillary networks—a smart way to get things connected. http://www.ericsson.com/res/thecompany/docs/publications/ericsson_review/2014/ericapillary-networks.pdf, 2014.
- ² BALLARDIE, T., FRANCIS, P., AND CROWCROFT, J. Core based trees (CBT). In *ACM SIGCOMM Computer Communication Review* (1993), vol. 23, ACM, pp. 85–95.
- ³ BONOMI, F., MILITO, R., ZHU, J., AND ADDEPALLI, S. Fog computing and its role in the internet of things. In *Proceedings of the first edition of the MCC workshop on Mobile cloud computing* (2012), ACM, pp. 13–16.
- ⁴ DAVIS, A., PARIKH, J., AND WEIHL, W. E. Edgecomputing: extending enterprise applications to the edge of the internet. In *Proceedings of the 13th international World Wide Web conference on Alternate track papers & posters* (2004), ACM, pp. 180–187.
- ⁵ GUPTA, T., SINGH, R. P., AND MAHAJAN, A. P. J. J. R. Bolt: Data management for connected homes. In *Proceedings of the 11th USENIX Symposium on Networked Systems Design and Implementation (NSDI 14)* (2014), pp. 243–256.
- ⁶ MORRIS, R., KOHLER, E., JANNOTTI, J., AND KAASHOEK, M. F. The click modular router. In *ACM SIGOPS Operating Systems Review* (1999), vol. 33, ACM, pp. 217–231.
- ⁷ RHEA, S. C., EATON, P. R., GEELS, D., WEATHERSPOON, H., ZHAO, B. Y., AND KUBIATOWICZ, J. Pond: The oceanstore prototype.
- ⁸ SATYANARAYANAN, M., BAHL, P., CACERES, R., AND DAVIES, N. The case for vm-based cloudlets in mobile computing. *Pervasive Computing, IEEE* 8, 4 (2009), 14–23.
- ⁹ SESKAR, I., NAGARAJA, K., NELSON, S., AND RAYCHAUDHURI, D. Mobilityfirst future internet architecture project. In *Proceedings of the 7th Asian Internet Engineering Conference* (2011), ACM, pp. 1–3.
- ¹⁰ TAMASSIA, R. Authenticated data structures. In *Algorithms-ESA 2003*. Springer, 2003, pp. 2–5.
- ¹¹ ZHANG, L., AFANASYEV, A., BURKE, J., JACOBSON, V., CROWLEY, P., PAPADOPOULOS, C., WANG, L., ZHANG, B., ET AL. Named Data Networking. *ACM SIGCOMM Computer Communication Review* 44, 3 (2014), 66–73.
- ¹² ZHAO, B. Y., KUBIATOWICZ, J., JOSEPH, A. D., ET AL. Tapestry: An infrastructure for fault-tolerant wide-area location and routing.