

Toward a Name-Based, Data-Centric Platform for Scientific Data

Lan Wang
University of Memphis
lanwang@memphis.edu

Lixia Zhang
UCLA
lixia@cs.ucla.edu

I. INTRODUCTION

Large-scale scientific measurements and experiments produce huge datasets that help advance science and improve human life [1]. For example, the sensors in our wearables, hand-held devices, and environments have been producing a massive amount of data that enables researchers to investigate a wide range of health and wellness issues. Although network transmission speed will continue to increase, this increase will likely lag behind the increase in the data volume. In addition, the number of data users is rapidly increasing, placing a high demand on the network for data distribution. In-situ computation can minimize unnecessary data transfers, but it requires identifying and utilizing computing resources close to the data sources for data processing. Moreover, as data may be processed through many stages at various distributed computing nodes, data users must be able to verify the integrity of computation results (i.e., data provenance). Finally, data privacy is an increasing concern in scientific research, especially for data involving human health and behavior. Unfortunately, it has been challenging to build effective solutions to meet the above efficiency and security requirements.

Several recent efforts have used Named Data Networking (NDN) [2] to support high volume transfer of climate data ([3], [4]) and high energy physics data [5]. NDN gives each piece of application data a unique, semantically meaningful name, and bind the name and data by a cryptographic signature (and encrypt it as needed). The named, secured data stays the same in the network, computing devices, and storage, which provides a strong foundation for supporting distributed computing, data provenance, and fine-grained access control.

II. FOUNDATION: NAMING AND SECURING DATA DIRECTLY

Scientific data is often produced and stored at facilities far away from the data users who may work in different organizations across the country, or even the globe. Currently there is no easy way to identify the relevant data in big data files, so users interested in specific pieces of data typically have to transfer the entire dataset, which leads to long delay, unnecessary bandwidth consumption, and server overload. This problem stems from the *lack of appropriate data naming* to facilitate efficient data fetching and processing, let alone data provenance and integrity checking. NDN addresses these challenges using hierarchically structured and semantically meaningful data names that can help with navigating a dataset and identifying potentially useful data. An NDN name can

represent a collection of one or more data items sharing the same name prefix. For example, `/org/md2k/mOral20` represents the dataset from an oral health study conducted by MD2K. It can be further divided into sub name prefixes (e.g., `/org/md2k/mOral20/gyro` and `/org/md2k/mOral20/accelerometer`), each representing a subset of the dataset (gyro or accelerometer data). These names can be further extended to identify individual pieces of data produced by a study participant at a specific time, e.g., `/org/md2k/mOral20/gyro/{user}/{timestamp}`. This means subsets of the data can be retrieved at any desired granularity.

Moreover, each data producer digitally signs its data to bind the data name to the data content. As such, data authenticity can be verified by anyone and each piece of data can be cached by, and served from any device, making data distribution more robust, efficient, and scalable. For example, when a user in US retrieves data from Australia, the data can be cached inside the US network. When other US users request the same data, the network can serve them the cached data, rather than fetching it from Australia again, and each user can verify the data's integrity and authenticity using its signature and associated trust model. *Note that NDN uses application-specific trust models that are more decentralized and robust than the global certificate infrastructure in use today.*

III. DISTRIBUTED COMPUTING OVER NDN

Król et. al. recently proposed “Compute First Networking” [6], an NDN-based distributed computing framework. They point out that today's distributed computing platforms rely on a complex set of centralized schedulers, DNS-based name translation, stateful load balancers, and heavy-weight transport protocols. An NDN-based system significantly simplifies the solution by employing *name-based network forwarding* and *multi-party data synchronization protocols (Sync)* [7]. Utilizing data names in network forwarding removes the need for name-to-address mapping and handles load balancing automatically. Unlike TCP/IP transport protocols that focus on point-to-point data delivery, Sync leverages NDN's unique and secured binding between name and content to achieve reliable data delivery among a *group of distributed nodes*. For example, if a group of compute nodes are tasked to detect biomarkers in the mOral20 data, they can form a sync group with the data producers (sensors) using the name prefix `/org/md2k/mOral20` and compare their set of data names to retrieve any missing data of interest using the data names. This approach is a natural fit for the many-to-many communication pattern that is common in a distributed computing scenario.

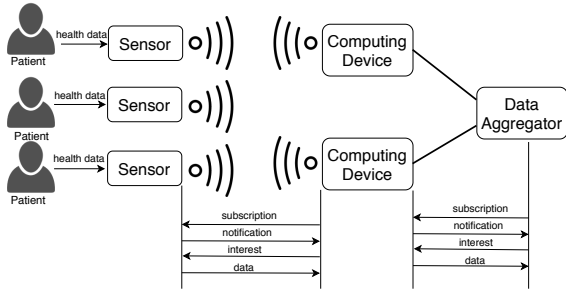


Fig. 1. Distributed Computing of Health Data from Sensors through Pub-Sub over NDN

Nichols [8] proposed a publication-subscribe API based on the NDN Sync primitive and used this pub-sub API to build a network measurement application. We believe that this pub-sub API can also be used in distributed computing to match computation requests and computing nodes. For example, as health data is collected from sensors on patients, the sensors can send computation requests (e.g., to detect biomarkers from the data) and any nearby computing devices can subscribe to such requests (Figure 1). Similarly, a data aggregator in the hospital can subscribe to the computed results.

IV. DATA PROVENANCE

Scientific research not only collects raw data from the environment and experiments, but also generates new data based on the raw data. As data from various sources goes through many iterations of aggregation, computation, and distillation, it becomes increasingly difficult to keep track of how a certain set of data is produced. In NDN, data is immutable – any computation on a piece of data produces a new piece of data which will have a new name. The application naming scheme can associate the original data name(s), computing node’s name, and other information with the new data name. Additional information can be published as meta data and given a name related to the primary data’s name. For example, if eating episodes are inferred from the mOral20 study’s gyroscope and accelerometer data (*/org/md2k/mOral20/gyro/(user)/(timestamp)* and */org/md2k/mOral20/accelerometer/(user)/(timestamp)*), then the eating episodes data’s can be named */org/md2k/mOral20/eating/(user)/(starttimestamp)-(endtimestamp)/(compute-node)*. The meta data can be named */org/md2k/mOral20/eating-metadata/(user)/(starttimestamp)-(endtimestamp)/(compute-node)* and contain the list of the gyroscope and accelerometer data names, computation time, inference algorithm, algorithm parameters, etc.

The application naming scheme, data names and associated meta data can be used to trace the series of computation and input data that led to a piece of data. In addition, the data signing key indicates who produced the data, and NDN can use trust schemas to automatically verify whether the owner of the key is authorized to produce the data [9].

V. ACCESS CONTROL

NDN names can facilitate both the specification and automated enforcement of access control policies. An access

control policy needs to identify *who are given access to what dataset(s) with what restrictions*. The data users can be named based on their organizations, e.g., */edu/memphis/lanwang*. Datasets can be identified using their names (e.g., */org/md2k/mOral20/gyro*), and restrictions can be specified using names of the data attributes and their ranges. Once the access control policy is defined, its enforcement is also *data-centric* without relying on SSL sessions or firewalls. Every piece of data is encrypted with a content key (C-KEY) that changes depending on the data and access granularity, and data access is controlled by encrypting and publishing the C-KEY for only those users authorized to access the data. We call this approach Name-based Access Control (NAC) [10]. The difference between NAC and traditional access control schemes is two-fold: (a) NAC ensures that data is encrypted both in transit and in storage, achieving true end-to-end protection; and (b) it can be applied to data at every granularity, from an entire data repository to a single data packet, following the hierarchical name structure.

VI. SUMMARY

Based on the NDN architecture that features *hierarchical naming structure, data centric security, and name-based data distribution*, the NDN project has generated a rich set of solutions with open source codebase that supports data navigation and discovery, distributed computing, data provenance, and access control, all of which are important for data-driven research. We are interested in collaborating with domain scientists, database researchers, and scientific software developers to build a name-based data-centric software platform for large-scale scientific research.

REFERENCES

- [1] R. Gerber, J. Hack, K. Riley, K. Antypas, R. Coffey, E. Dart, T. Straatsma, J. Wells, D. Bard, S. Dosanjh *et al.*, “Crosscut report: Exascale requirements reviews,” Oak Ridge National Lab.(ORNL), Oak Ridge, TN (United States); Argonne, Tech. Rep., 2018.
- [2] L. Zhang, A. Afanasyev, J. Burke, V. Jacobson, K. Claffy, P. Crowley, C. Papadopoulos, L. Wang, and B. Zhang, “Named Data Networking,” *ACM SIGCOMM Computer Communication Review*, July 2014.
- [3] C. Olschanowsky, S. Shannigrahi, and C. Papadopoulos, “Supporting Climate Research using Named Data Networking,” in *IEEE LANMAN*, 2015.
- [4] S. Shannigrahi, C. Fan, and C. Papadopoulos, “Request Aggregation, Caching, and Forwarding Strategies for Improving Large Climate Data Distribution with NDN: A Case Study,” in *In Proceedings of ACM ICN*, September 2017.
- [5] S. Shannigrahi, A. Barczyk, C. Papadopoulos, A. Sim, I. Monga, H. Newman, J. Wu, and E. Yeh, “Named Data Networking in Climate Research and HEP Applications,” in *21st International Conference on Computing in High Energy and Nuclear Physics (CHEP2015)*, 2015.
- [6] M. Król, S. Mastorakis, D. Oran, and D. Kutscher, “Compute first networking: Distributed computing meets icn,” in *Proceedings of the 6th ACM Conference on Information-Centric Networking*, 2019, pp. 67–77.
- [7] T. Li, W. Shang, A. Afanasyev, L. Wang, and L. Zhang, “A Brief Introduction to NDN Dataset Synchronization (NDN Sync),” in *IEEE MILCOM*, 2018.
- [8] K. Nichols, “Lessons learned building a secure network measurement framework using basic ndn,” in *Proceedings of the 6th ACM Conference on Information-Centric Networking*, 2019, pp. 112–122.
- [9] Y. Yu, A. Afanasyev, D. Clark, K. Claffy, V. Jacobson, and L. Zhang, “Schematizing and automating trust in Named Data Networking,” in *ACM ICN*, September 2015.
- [10] Z. Zhang, Y. Yu, S. K. Ramani, A. Afanasyev, and L. Zhang, “NAC: Automating access control via Named Data,” in *IEEE MILCOM*, 2018.