# HiFi: A Unified Architecture for High Fan-in Systems

## (System Demonstration)

Owen Cooper*, Anil Edakkunni*, Michael J. Franklin*, Wei Hong[+], Shawn R. Jeffery*,
Sailesh Krishnamurthy*, Fredrick Reiss*, Shariq Rizvi*, and Eugene Wu*

*EECS Dept., UC Berkeley      [+]Intel Research Berkeley

## Abstract

Advances in data acquisition and sensor technologies are leading towards the development of "High Fan-in" architectures: widely distributed systems whose edges consist of numerous receptors such as sensor networks and RFID readers and whose interior nodes consist of traditional host computers organized using the principle of successive aggregation. Such architectures pose significant new data management challenges. The *HiFi* system, under development at UC Berkeley, is aimed at addressing these challenges. We demonstrate an initial prototype of *HiFi* that uses data stream query processing to acquire, filter, and aggregate data from multiple devices including sensor motes, RFID readers, and low power gateways organized as a High Fan-in system.

## 1. Introduction

Emerging wireless sensor networks and RFID technologies are on a fast track to widespread deployment in applications such as environmental monitoring, asset tracking, telemetry-based remote monitoring, and real-time supply chain management. Driven by both technological and market forces, sensor and RFID-based deployments raise the promise of taking computing from its current, user-driven mode, to one of direct and continuous interaction with the physical world.

**Proceedings of the 30th VLDB Conference,
Toronto. Canada. 2004**

### 1.1 High Fan-in Systems

In many cases, sensors and readers will serve as the receptors at the edges of widely distributed systems. For example, in a supply chain management deployment, collections of sensors and RFID readers on individual store shelves (in a retail scenario) or dock doors (in a warehouse/manufacturing scenario) continuously collect readings. These readings include "beeps" from low-function passive RFID tags (indicating the presence of particular tagged objects, such as cases of goods), as well as more content-rich information from smart sensors and higher-function tags (such as temperature readings, shipping histories, etc.).

These "edge" devices produce data that will be aggregated locally with data from other nearby devices. That data will be further aggregated within a larger area, and so on. This arrangement results in a distinctive bowtie topology we refer to as a High Fan-In system (see Figure 1). A sophisticated system such as one supporting a nation-wide supply chain application may consist of vast numbers of widely dispersed receptor devices (i.e., many thousands or more depending on the technology) and many levels of successively wider-scoped aggregation and storage. Such systems will comprise a vast array of heterogeneous resources, including inexpensive tags, wired and wireless sensing devices, low-power compute nodes and PDAs, and computers ranging from laptops to the largest mainframes and clusters.
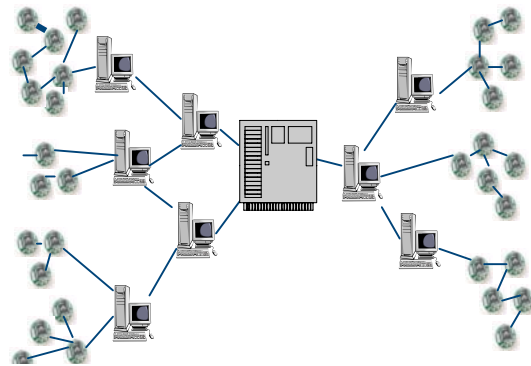


**Figure 1 - A High Fan-in Architecture**

## 1.2 Query Processing: A Unifying Framework

Traditionally, sensor-based information systems have been deployed using a piecemeal approach — a sensor-specific programming environment is used to task the edge receptors, a separate transport or information bus is used to route the sensor readings, and a database system or other data manager is used to collect and process the sensor readings. As a result, sensor application deployments have tended to be costly, difficult, and inflexible.

In contrast, our work is based on the notion that database techniques in general and data stream query processing in particular have a role to play *at all levels* of large-scale High Fan-in systems. In fact, we propose to use stream query processing and views as a *unifying framework* for data access across an entire High Fan-in environment. Thus, while it is certainly the case that the various types of sensor devices and computing platforms present in such systems all have their own unique characteristics and idiosyncrasies, our design uses stream query processing as the glue that binds these disparate pieces together to become a highly-functional application deployment platform.

Building on earlier work at Berkeley on both adaptive data stream processing (TelegraphCQ [CCDF+03]) and sensor network databases (TinyDB [MFHH03]) we have developed a new multi-level system for continuous and historical data processing in large-scale, sensor-rich applications. Our system, called HiFi, consists of host-centric processing components and device-centric processing components that can interoperate in a seamless manner. As such, HiFi runs across a gamut of platforms, ranging from battery powered wireless sensor motes and RFID readers, to mid-tier Linux-based sensor platforms and PDAs, to plugged-in, diskful servers.

## 2. HiFi Challenges and Architecture

While building on the growing body of work in the areas of data stream processing, sensor network databases, and data integration, the design of HiFi also addresses a number of challenges that arise from the unique properties of High Fan-in architectures and the applications they support.

### 2.1 The Challenges of Scale

From a data management perspective, the most challenging new aspect that High Fan-in systems bring to the table is the tremendous range they span in terms of three key characteristics: time, space, and resources.

**Time** – Timescales of interest in a High Fan-in system can range from seconds or less at the edges, to years in the interior of the system. At the edges of a High Fan-in system are the receptor devices that repeatedly measure some aspect of the physical world (or perhaps the virtual world, say in the case of a network monitor or other logical sensor). These devices are typically concerned with fairly short time scales, perhaps on the order of seconds or less. As one moves away from the edges, the timescales of interest increase.

For example, in a retail RFID scenario, individual readers on shelves may read several times a second, while the manager of a store may be concerned with how sales of particular items are going over the course of a morning, and planners at regional and corporate centers may be more concerned with longer-term sales trends over a season or several seasons.

**Space** - As with time, the area of geographic interest grows significantly as one moves from the edges of a High Fan-in system to the interior. Again using the retail RFID scenario, individual readers are concerned with a space of a few square meters, aggregation points within the store would be concerned with entire departments or perhaps the store as a whole, and regional and national centers are concerned with those much larger geographical areas.

**Resources** – Finally, the range of computing resources available at various levels of a High Fan-in system also vary dramatically, from small, cheap sensor motes (e.g., Berkeley Motes, see Figure 2) on the edges up to the largest mainframes in the interior of the system. Communication resources also can range from low power, lossy radios at the edges, to dedicated high-speed fiber in the interior.



**Figure 2 - Berkeley Mica Wireless Sensor Mote**

The key architectural principle on which the design of HiFi rests is the use of declarative query processing to provide *uniform* data access across all of these various scales. As with other data-intensive environments, the idea is to use the declarative approach to shield application developers from the complexities of the

underlying platform, while enabling the system to optimize and efficiently execute data access and processing operations.

## 2.2 Cascading Stream Processing

HiFi is based on a notion of cascading stream processing, in which data streams collected at the very edges of the network are continually filtered, refined, aggregated and brought in towards the interior of the network. The multi-platform nature of the system enables seamless transitions between sensor-like devices, mid-tier single-board sensor aggregators, and host computers.

Each RFID reader or sensor network access point in a HiFi network is a data stream generation machine at the edge of a large (potentially global) information system. While the volume of data from any single stream is likely to be modest, the low cost and eventual ubiquity of such devices leads to a torrent of data in aggregate.

Furthermore, because these devices are embedded in the physical world, they are geographically proximate to the entities or spaces they are intended to be monitoring. As such, the processing of the data streaming from these devices will exhibit highly-localized patterns at many granularities. For example, an RFID reader on a store shelf sampling several times a second will constantly and repeatedly detect the presence of items on that shelf. While at one level, these "beeps" are indeed a data stream, in general the continued presence of an item on the shelf is not data of interest beyond the scope of that shelf. Rather, it is only when an item is removed or a new item appears that an interesting event worthy of propagation can be said to occur.

Such concerns argue for the ability to place stream processing at or close to the edges of these receptor-based systems, in order to perform highly-local tasks such as data cleaning, filtering, and simple event detection. Another argument for pushing stream query processing out towards the edges of the network is to reduce the communication requirements for battery-powered wireless devices such as Berkeley Motes, as has been demonstrated in our TinyDB work and other sensor network database projects.

On the other hand, much stream processing is likely to be best performed on mid-tier devices such as the Intel Stargate single-board computer (see Figure 3), which is a low-power Linux machine that can run on Li Ion batteries or can be plugged into the wall, and can use 802.11 communication (compared to the lower range/lower power radios used on sensor motes). Processing tasks that involve the correlation of multiple streams and/or the application of more sophisticated filtering and business rules are good candidates for such devices.

Still other stream processing will best be done on host computers. As you move in towards the interior of the network, three factors combine to change the nature of

device requirements. First, these nodes serve as aggregation points for ever larger geographic regions, resulting in an increase in the number of streams being monitored and the aggregate data volume to be processed. Secondly, many applications will need to archive streams and provide access to that past data. Thus, diskful systems will be required. Finally, the availability of large volumes of both live and historical stream data will make such systems magnets for queries from throughout the network.



**Figure 3 - Stargate Mid-tier Processing Node**

The hierarchical nature of a High Fan-in architecture leads naturally to the use of stream query processing for *aggregation*. That is, as data streams flow from the edges towards the interior they can be combined and aggregated in order to produce summaries and reduce data volume. Indeed, aggregation is one of the major uses of stream query processing in HiFi. There are, however, a number of other important tasks for which we believe stream query processing is particularly well suited. These include:

- **Data Cleaning** – Sensors and RFID readers are notoriously noisy devices, and dealing with the poor quality of data they produce is one of the main challenges in a High Fan-in (or any sensor-based) system. We believe that declarative queries can be used to specify cleaning functionality for any single device as well as across groups of devices.

- **Event Monitoring** – One of the main functions of a High Fan-in system is to continuously monitor the edge environment, and to send alerts when events of interest are detected. These events may vary significantly in terms of the timescales and geographic areas over which they are detected. While many commercial systems

1359

have developed their own "event languages", we believe that a stream query language (perhaps suitably extended) is the right substrate for such functionality.

- **Stream Correlation** – A further advantage of a query-based approach is that it natively supports the ability to compare and correlate data from multiple streams. Such streams may be homogeneous, as in the case of comparing temperature readings from a group of identical sensors, or heterogeneous, as in the case of combining temperature readings with RFID "beeps".

- **Outlier Detection** – Another form of data reduction provided by a HiFi System is outlier detection. In many monitoring systems, expected events are of less immediate interest than anomalies. Queries can be used to detect and propagate various types of outliers in a streaming environment.

As the above (partial) list indicates, we believe that stream query processing can serve as the basis of a wide range of High Fan-in functionality in a uniform manner. This is a significant departure from the *ad hoc* way in which such systems are currently being constructed. Our initial HiFi prototype has been developed as a proof-of-concept platform, in order to test the viability of our query-based approach.

## 3. Overview of the Demo

We have built an initial version of HiFi using the TelegraphCQ stream query processor and the TinyDB sensor database system. The goal of this prototype is to examine the feasibility of the uniform query processing model and to derive a better understanding of the core components required for building High Fan-in systems.

In this demo, we will show how to use HiFi to query and correlate data from multiple streaming sources. In particular we intend to implement a simple tracking application using passive RFID tags and sensor motes and showing the power of stream query processing for providing real-time analyses of correlated readings across multiple streams and device types and successive levels of aggregation. In addition to showing some gadgets such as RFID tags/readers and several classes of wireless sensor devices, we will also demonstrate some of the salient features of the architecture, including:

- Integration of multiple platforms including wireless sensor motes, passive RFID tags and

readers, mid-tier sensor aggregation points, and host computers.

- Unified and adaptive cross-platform query optimization in this heterogeneous environment.

- In-network query processing both in the interior and at the edges of the network.

- The use of stream queries and views for providing key functionality such as data cleaning, filtering, and event detection.

- Support for continuous multi-level aggregation.

We will show HiFi running on a heterogeneous network consisting of at least the following platforms:

1) Laptop PCs

2) Intel Stargate low-power single-board computer.

3) Berkeley sensor motes capable of sensing light, temperature, and sound.

4) RFID readers with passive, read-only RFID tags.

The demonstration will highlight the various features and underlying technology of the HiFi design and will emphasize the power and flexibility of a uniform stream query model for supporting applications over multiple classes of receptor devices and networks.

## REFERENCES

[CCDF+03] S. Chandrasekaran, O. Cooper, A. Deshpande, M. Franklin, J. Hellerstein, W. Hong, S. Krishnamurthy, S. Madden, V. Raman, F. Reiss, and M. Shah, "TelegraphCQ: Continuous Data Flow Processing for an Uncertain World", Proceedings of the 1st Conference on Innovative Data Systems Research (CIDR 2003), Asilomar, CA, January, 2003.

[MFHH03] S. Madden, M. Franklin, J. Hellerstein, and W. Hong, "The Design of an Acquisitional Query Processor for Sensor Networks", Proceedings of the ACM SIGMOD Int'l Conf. on Management of Data (SIGMOD 2003), San Diego, CA, June 2003, pp 491-502.