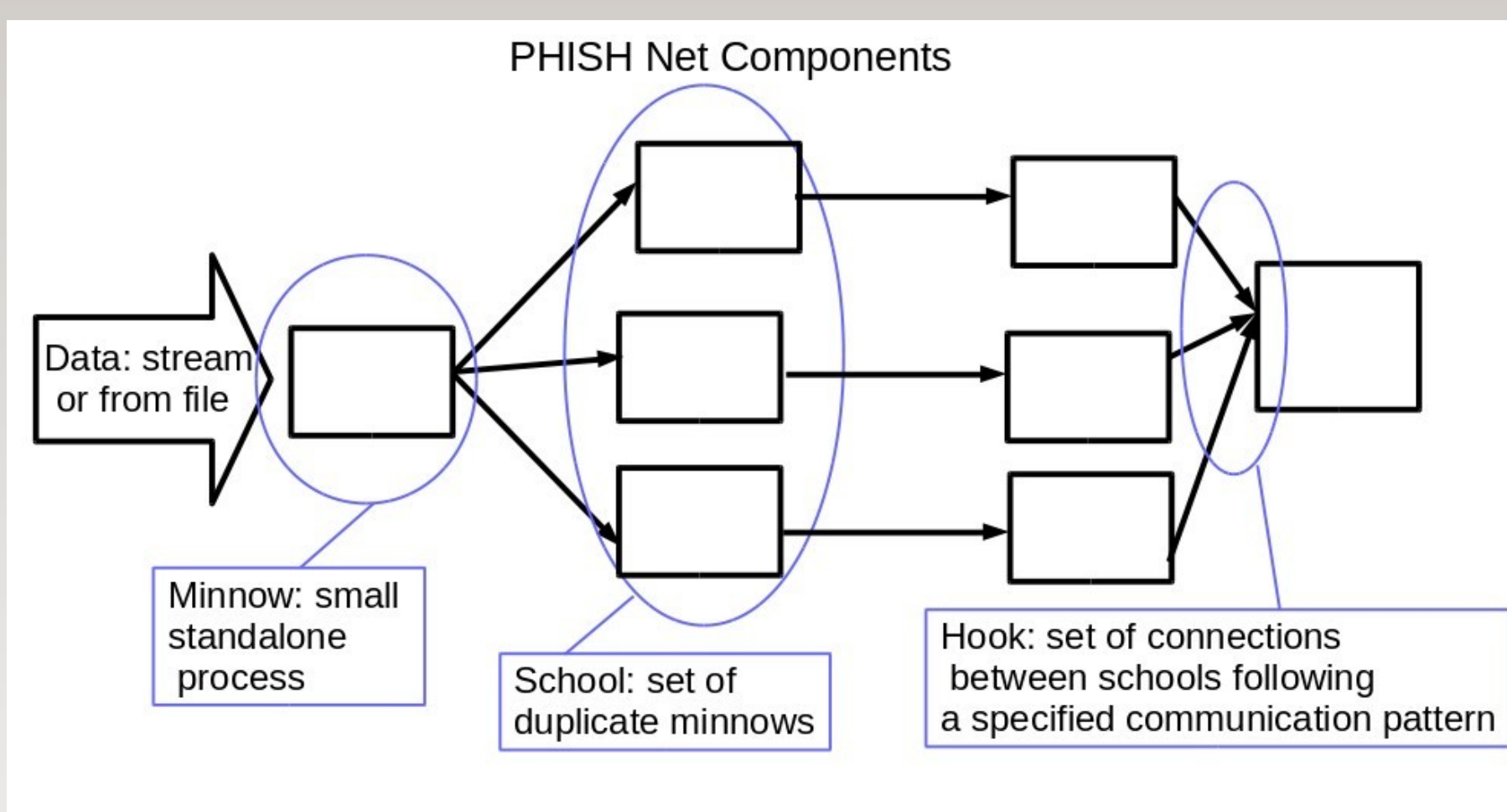# Data Structures for Parallel High Speed Streaming

Alexandra Porter, Arizona State University, Computer Science & Mathematics, Spring 2017
Manager: Jenn Troup, Mentors: Jon Berry, Org 1464, Cindy Phillips, Org 1400
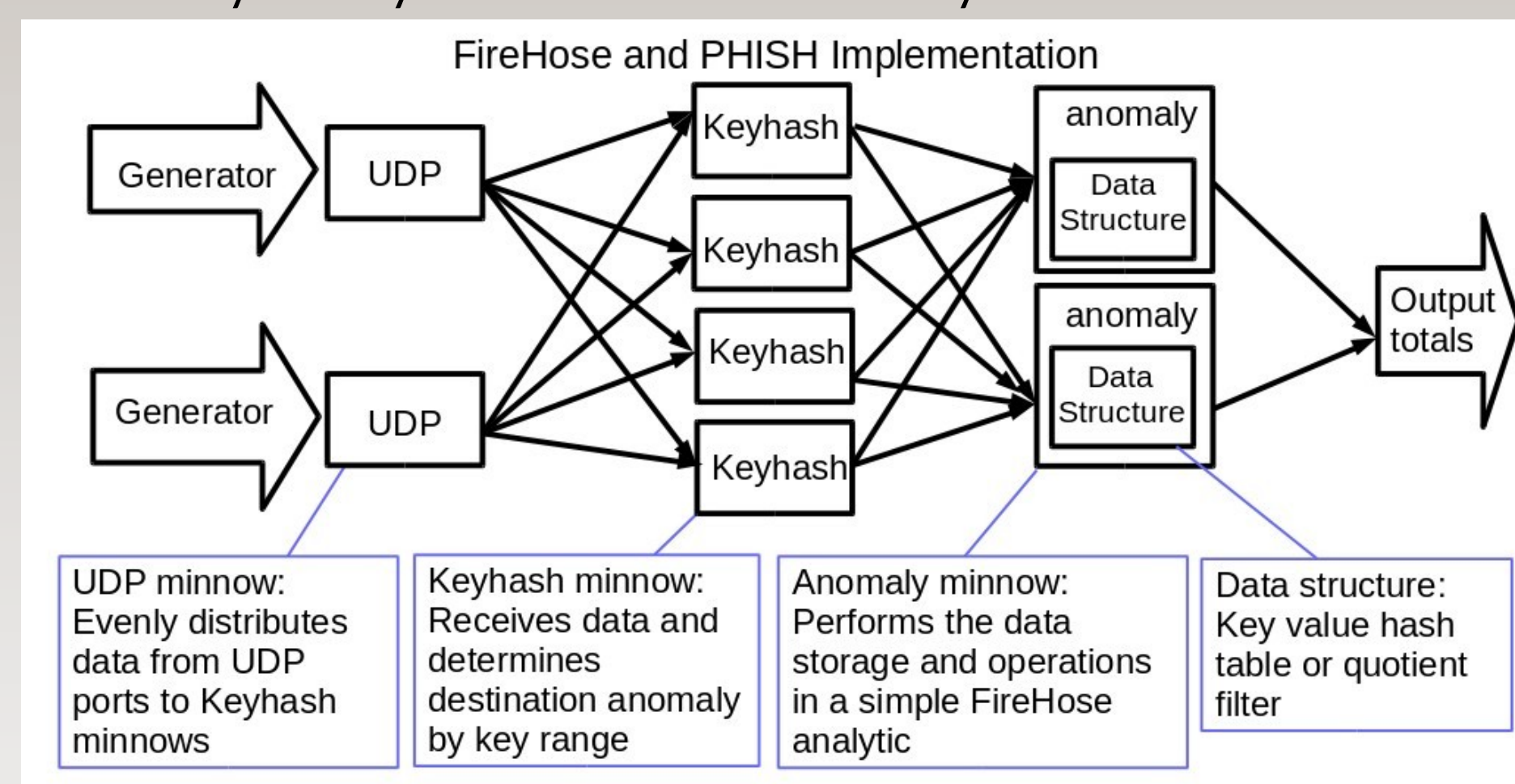Sandia National Laboratories/NM, U.S. Department of Energy, July 26, 2016

## PHISH: Parallel Harness for Informatic Stream Hashing

- Software Framework for high speed streaming
- Can run on multi-core and HPC machines

**PHISH Net Components**



- Data: stream or from file
- Minnow: small standalone process
- School: set of duplicate minnows
- Hook: set of connections between schools following a specified communication pattern

## FireHose Streaming Benchmark

- A set of generators and a set of analytics
- Generators produce streams of packets written to UDP ports
- Analytics read packets, store state derived from key-value stream
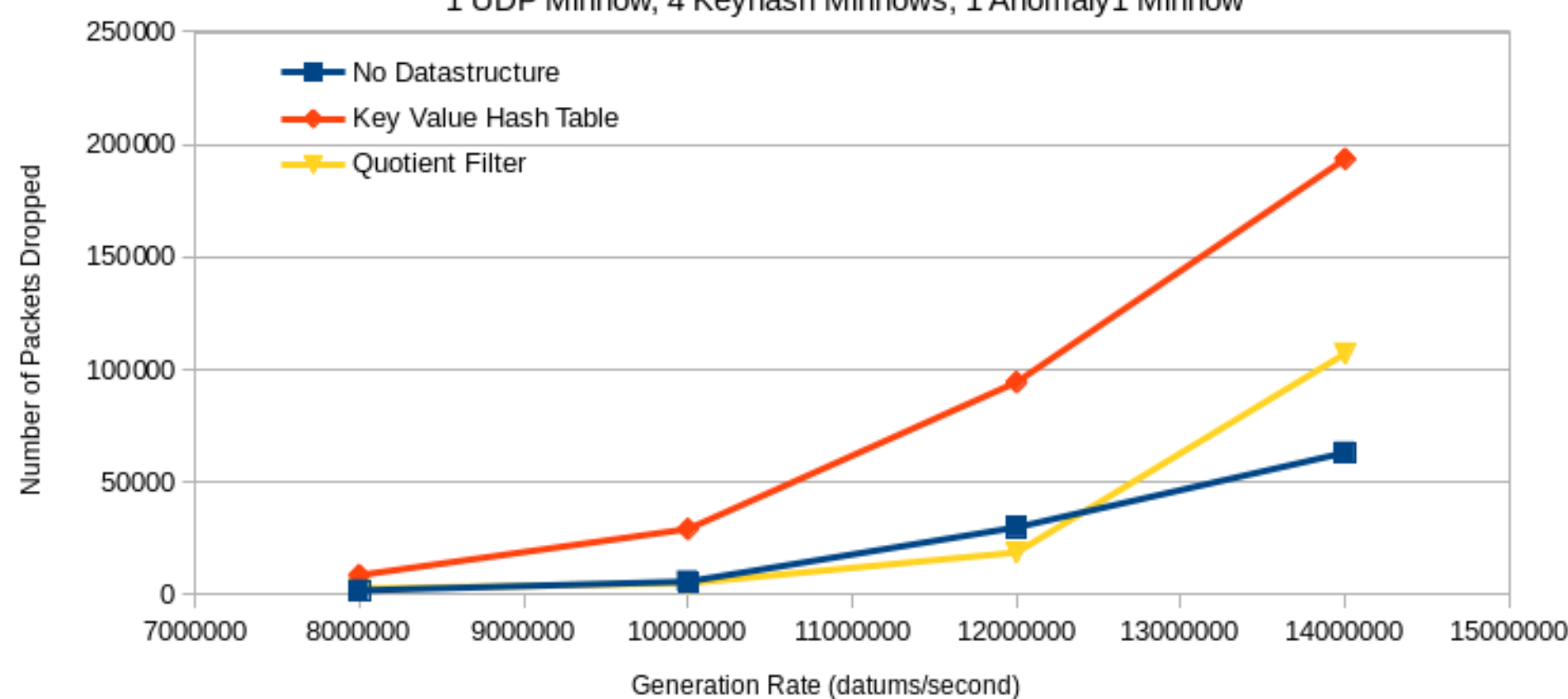- Anomaly1 analytic determines which keys are biased

**FireHose and PHISH Implementation**



- UDP minnow: Evenly distributes data from UDP ports to Keyhash minnows
- Keyhash minnow: Receives data and determines destination anomaly by key range
- Anomaly minnow: Performs the data storage and operations in a simple FireHose analytic
- Data structure: Key value hash table or quotient filter

## Research Challenges

- Configuring PHISH for optimal parallelism
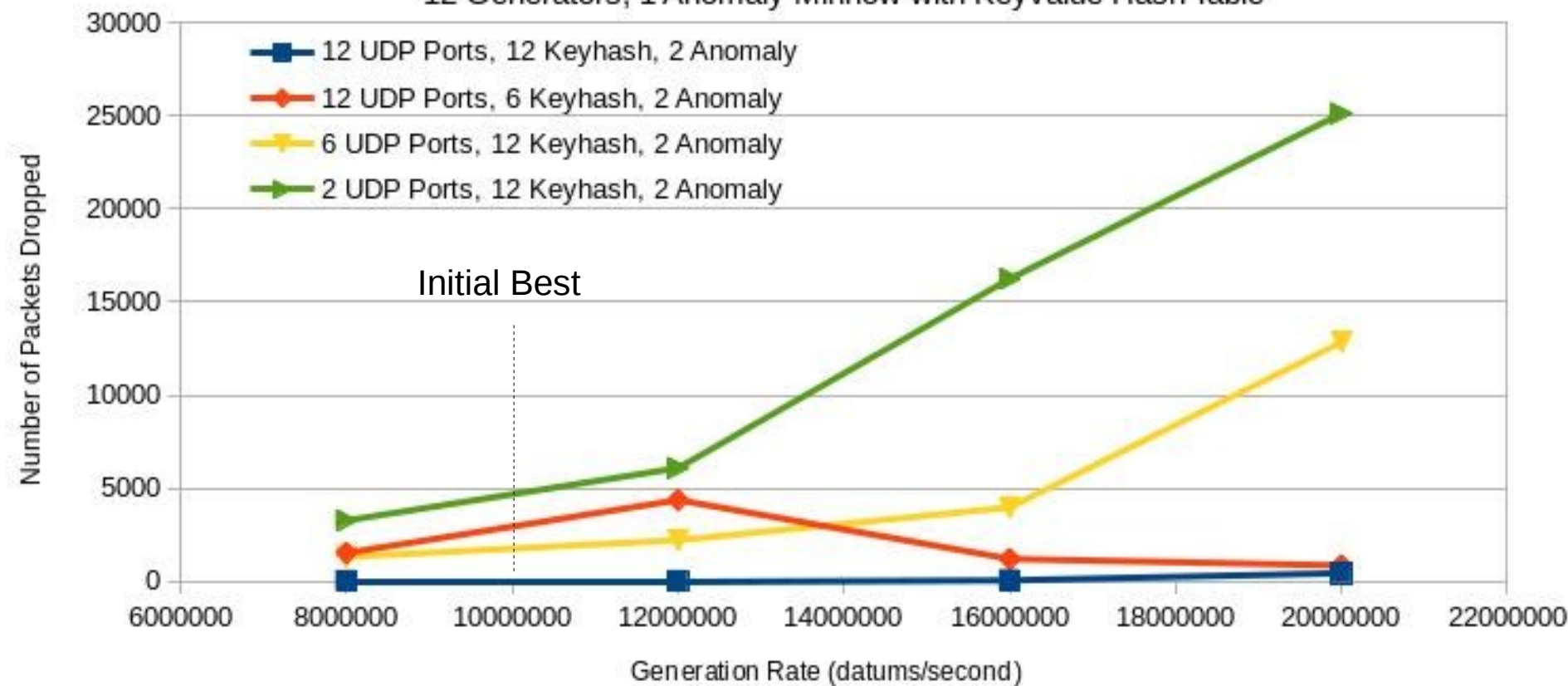- Optimizing data structures to maximize throughput

## Preliminary Results

- PHISH Anomaly1 analytic implementation handles up to 20 million key-value pairs per second on dual-processor Haswell
- Data structure choice crucial for further scaling



Comparison of Anomaly Data Structures
1 UDP Minnow, 4 Keyhash Minnows, 1 Anomaly1 Minnow



Comparison of PHISH Networks
12 Generators, 1 Anomaly Minnow with Keyvalue Hash Table

## Approximate Membership Data Structures

- Bloom Filters: fast membership testing but may return false positives
- Quotient Filters: allow counting and deletion of entries

**Simple Quotient Filter**

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| occupieds | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| runends | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 |
| remainders | a | b | c | d | e | | f | g | |

## Our Contributions to Quotient Filters

- Incorporate time stamps with constant number of bits
- Systematic expiration strategy to support advanced FireHose analytics

**Quotient Filter Aging Cycle**

| Time: | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Insertion time stamp: | 0 | 1 | 2 | 3 | 4 | 0 | 1 | 2 | 3 | 4 | 0 | 1 |
| | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| Deleting entries with stamp: | | | | | | 0 | 1 | 2 | 3 | 4 | 0 | 1 |
| | | | | | | 0 | 0 | 0 | 0 | 0 | 1 | 1 |

FireHose: firehose.sandia.gov/
PHISH: Streaming data analytics via message passing with application to graph algorithms, S. J. Plimpton and T. Shead, J Parallel and Distributed Computing, 74, 2687-2698 (2014).
www.sandia.gov/~sjplimp/phish.html