# Pivot Tracing
## Dynamic Causal Monitoring for Distributed Systems

**Jonathan Mace**, Ryan Roelke, Rodrigo Fonseca

Brown University

# Pivot Tracing
## Dynamic Causal Monitoring for Distributed Systems

Dynamically instrument live distributed systems

# Pivot Tracing
## Dynamic Causal Monitoring for Distributed Systems

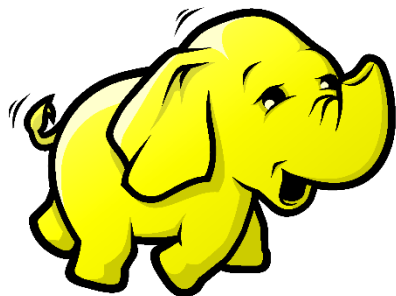Dynamically instrument live distributed systems

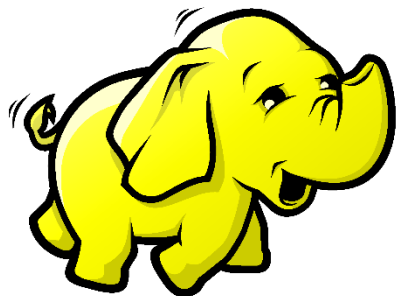Correlate and group events across components

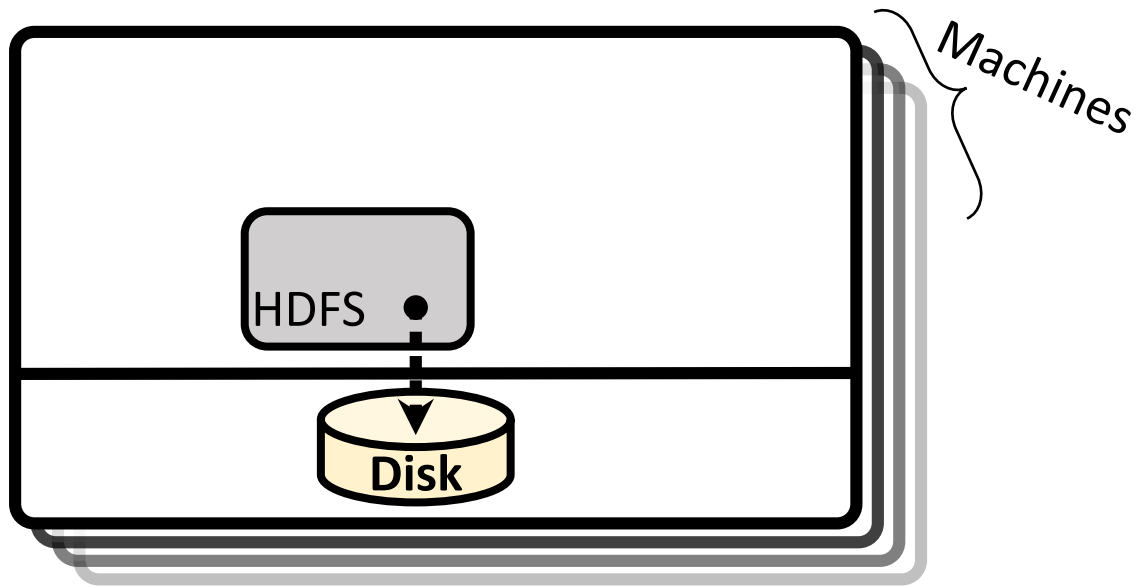# Pivot Tracing
## Dynamic Causal Monitoring for Distributed Systems

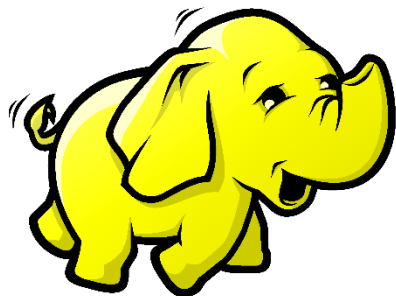Dynamically instrument live distributed systems
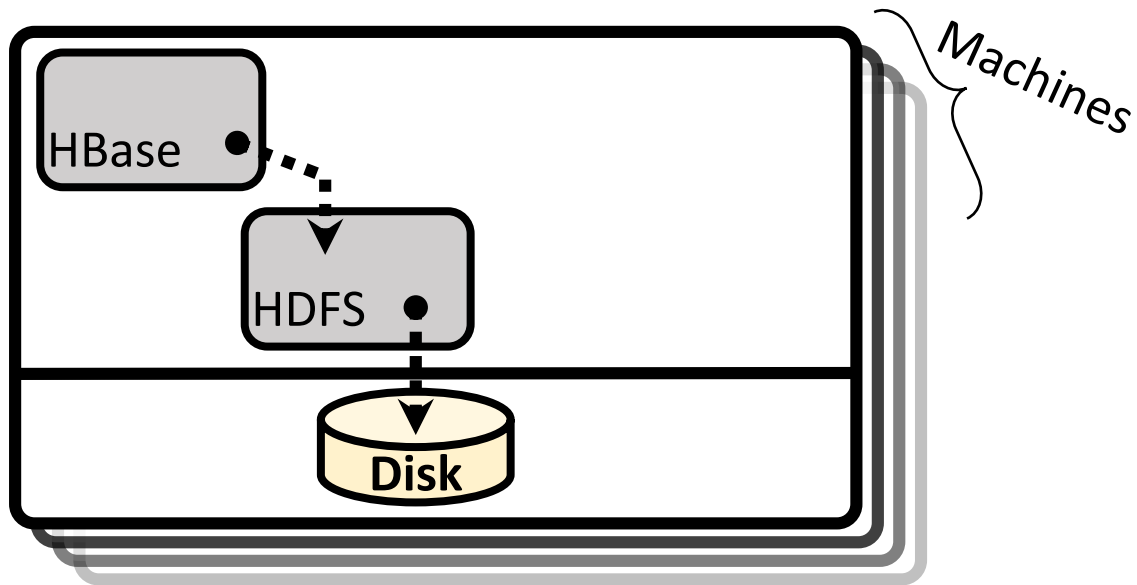
Correlate and group events across components

Hadoop Stack

Machines

HDFS •

Disk

Hadoop Stack

Machines

HBase

HDFS

Disk

Hadoop Stack

3

Machines

HBase
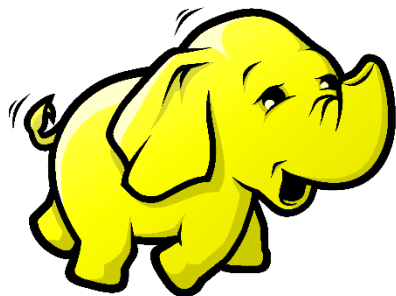
MapReduce
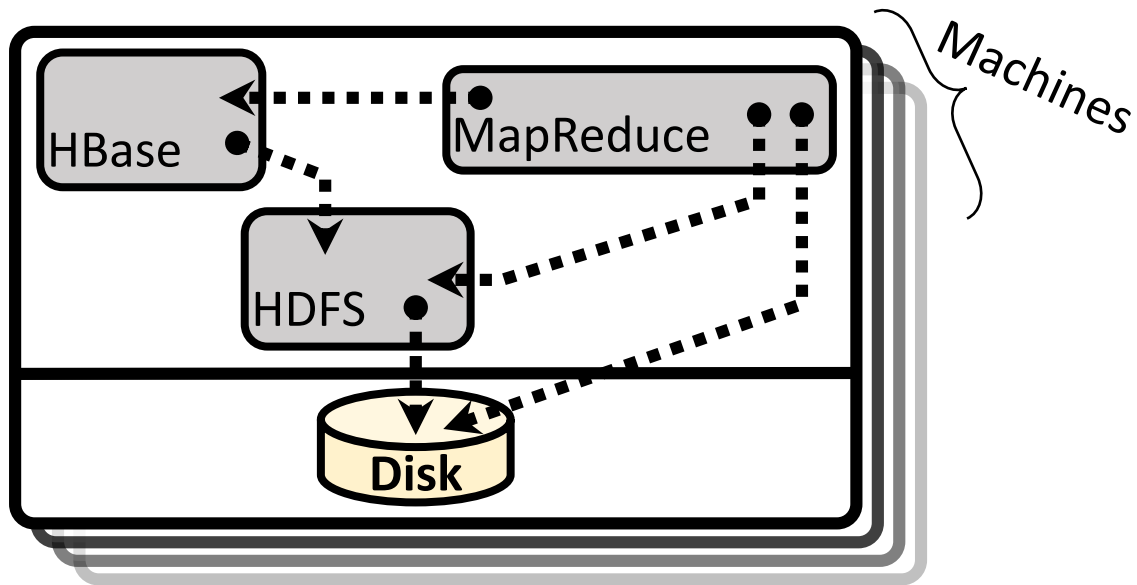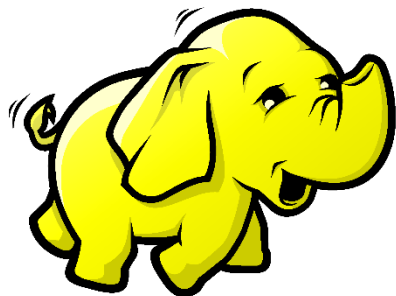
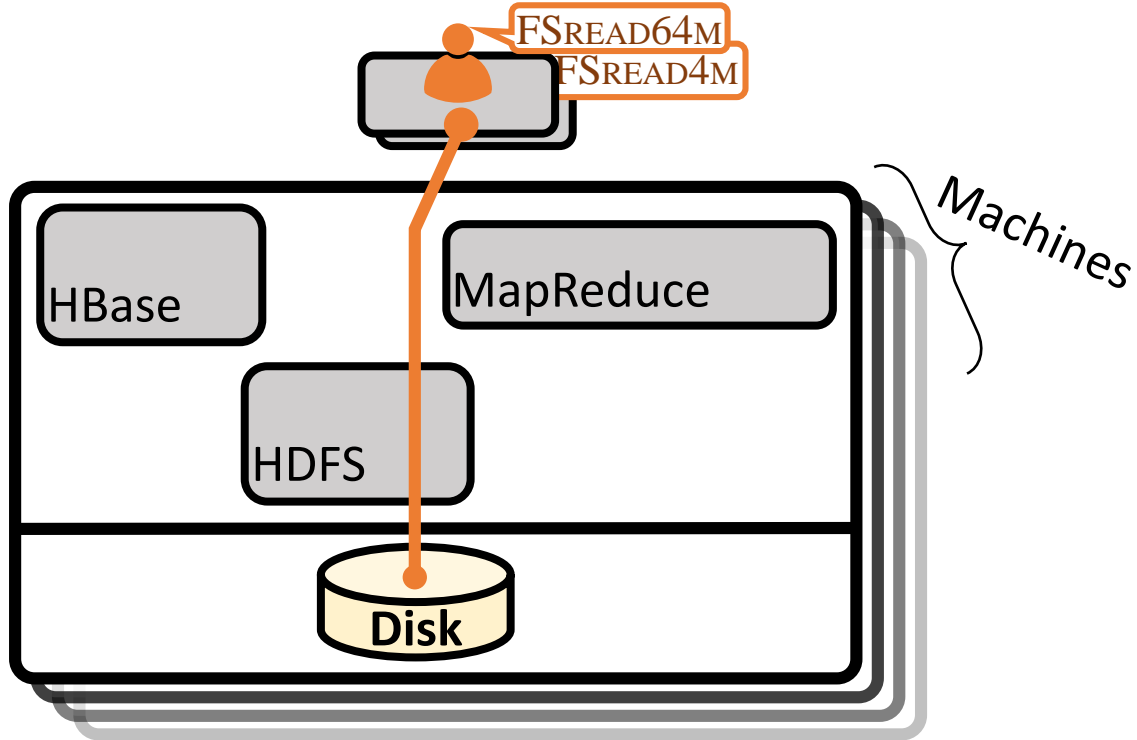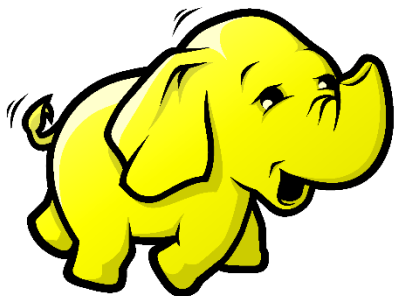HDFS

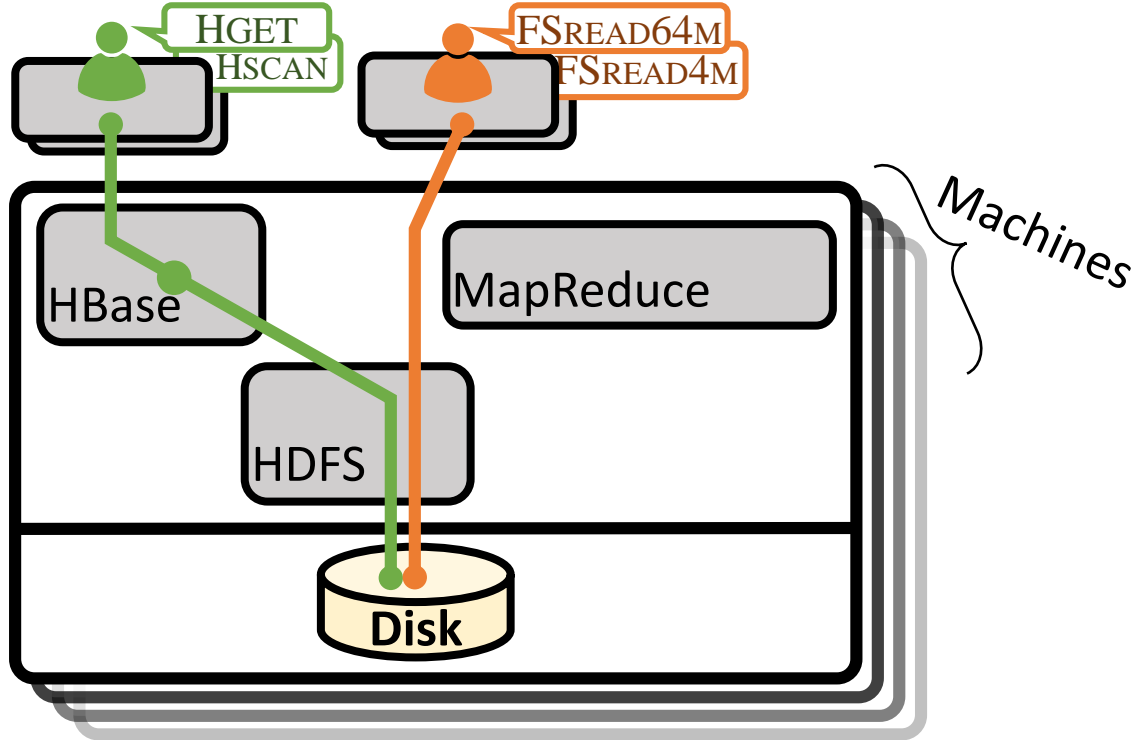Disk

Hadoop Stack

FSREAD64M
FSREAD4M

HBase

MapReduce

HDFS

Disk

Machines

Hadoop Stack

Hadoop Stack

Hadoop Stack

How is disk bandwidth being used?

How is disk bandwidth being used?

How is disk bandwidth being used?

# How is disk bandwidth being used?

How is disk bandwidth being used?

How is disk bandwidth being used?

How is disk bandwidth being used?

HGET
HSCAN

FSREAD64M
FSREAD4M

MRSORT100G
MRSORT10G

Machines

HBase

MapReduce

HDFS

IOStream

Disk

HDFS Throughput [MB/s]

DataNode
ReduceTask
MapTask
ShuffleHandler

200

150

100

50

0

0        5        10        15

Time [min]

9

HGET
HSCAN

FSREAD64M
FSREAD4M

MRSORT100G
MRSORT10G

Machines

HBase

MapReduce

HDFS

IOStream

Disk

MRSORT10G
MRSORT100G
FSREAD64M
FSREAD4M
SORT10G  (ShuffleHandler)
SORT10G  (ReduceTask)
SORT100G (MapTask)

HGET
HSCAN

HDFS Throughput [MB/s]

Time [min]

# Instrumentation is decided at development time

# Instrumentation is decided at development time

Probably not have enough info for your problem

Probably too much irrelevant info for your problem

# Instrumentation is decided at development time

Probably not have enough info for your problem

Probably too much irrelevant info for your problem

Should every user bear the cost of a feature?

# Instrumentation is decided at development time

Probably not have enough info for your problem

Probably too much irrelevant info for your problem

Should every user bear the cost of a feature?

HDFS-6292 Display HDFS per user and per group usage on webUI.

https://issues.apache.org/jira/browse/HDFS-6292

# Instrumentation is decided at development time

Probably not have enough info for your problem

Probably too much irrelevant info for your problem

Should every user bear the cost of a feature?

HDFS-6292 Display HDFS per user and per group usage on webUI.

https://issues.apache.org/jira/browse/HDFS-6292

# Dynamic dependencies

You often need to correlate information from different points in the system

# Dynamic dependencies

You often need to correlate information from different points in the system

Systems are designed to compose

Systems don't embed monitoring that relates to other services

# Dynamic dependencies

You often need to correlate information from different points in the system

Systems are designed to compose

Systems don't embed monitoring that relates to other services



Netflix "Death Star" Microservices Dependencies

@bruce_m_wong

11

# Pivot Tracing

# Pivot Tracing

You don't know the questions in advance

**Dynamic instrumentation**

**Fay (SOSP'11), Dtrace (ATC'04), …**

# Pivot Tracing

You don't know the questions in advance

**Dynamic instrumentation**

**Fay (SOSP'11), Dtrace (ATC'04), …**

You often need to correlate information from different points in the system

**Causal tracing**

**X-Trace (NSDI'07), Dapper (Google), Pip (NSDI'06), …**

# Pivot Tracing

# Pivot Tracing

Model system events as tuples in a streaming, distributed dataset

# Pivot Tracing

Model system events as tuples in a streaming, distributed dataset

Dynamically evaluate relational queries over this dataset

# Pivot Tracing

Model system events as tuples in a streaming, distributed dataset

Dynamically evaluate relational queries over this dataset

## Happened-before Join ( ⋈⃗ )
Join based on Lamport's *happened-before* relation

# Pivot Tracing
## Overview

HBase

MapReduce

DataNode
Metrics

HDFS

```
DataNodeMetrics.java
50    public class DataNodeMetrics {
        ...
266      public void incrBytesRead(int delta) {
267        ...
268      }
        ...
407  }
```

DataNode
Metrics

HBase

MapReduce

HDFS

```
DataNodeMetrics.java
50     public class DataNodeMetrics {
         ...
266      public void incrBytesRead(int delta) {
267          ...
268      }
         ...
407    }
```

DataNode Metrics

("DataNodeMetrics", delta=10, host="hop01", ...)

HBase

MapReduce

HDFS

DataNodeMetrics.java

```java
50    public class DataNodeMetrics {
         ...
266      public void incrBytesRead(int delta) {
267         ...
268      }
         ...
407   }
```

DataNode Metrics

("DataNodeMetrics", delta=10, host="hop01", …)

HBase

MapReduce

HDFS


HDFS Throughput [MB/s] vs Time [min]

Legend: Host A, Host B, Host C, Host D, Host E, Host F, Host G, Host H

```
From incr In DataNodeMetrics.incrBytesRead
   GroupBy incr.host
   Select incr.host, SUM(incr.delta)
```

```
DataNodeMetrics.java
50    public class DataNodeMetrics {
         ...
266      public void incrBytesRead(int delta) {
267          ...
268      }
         ...
407  }
```
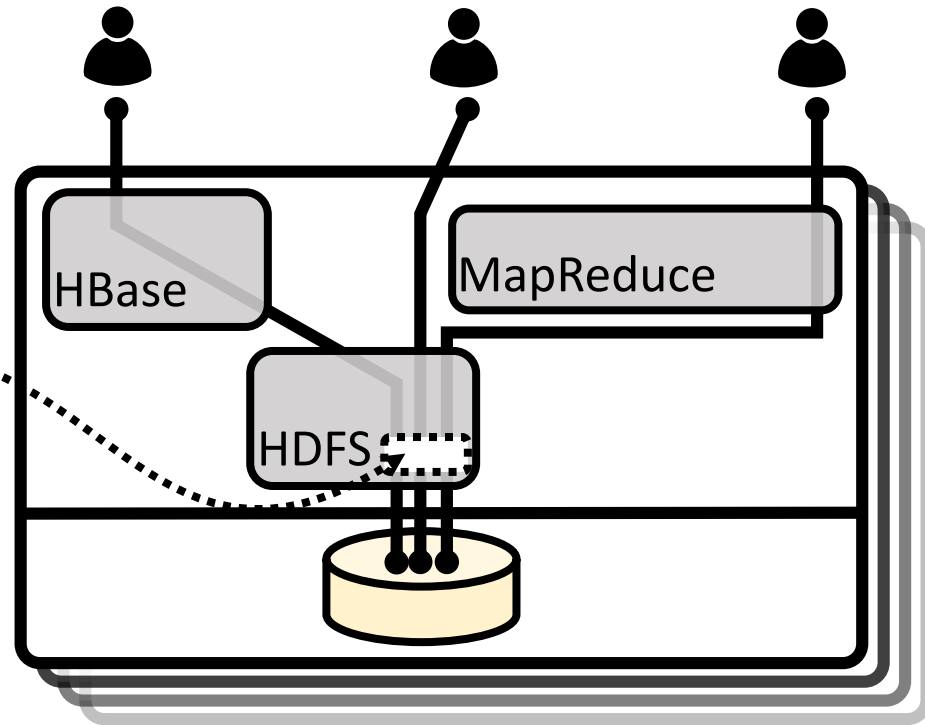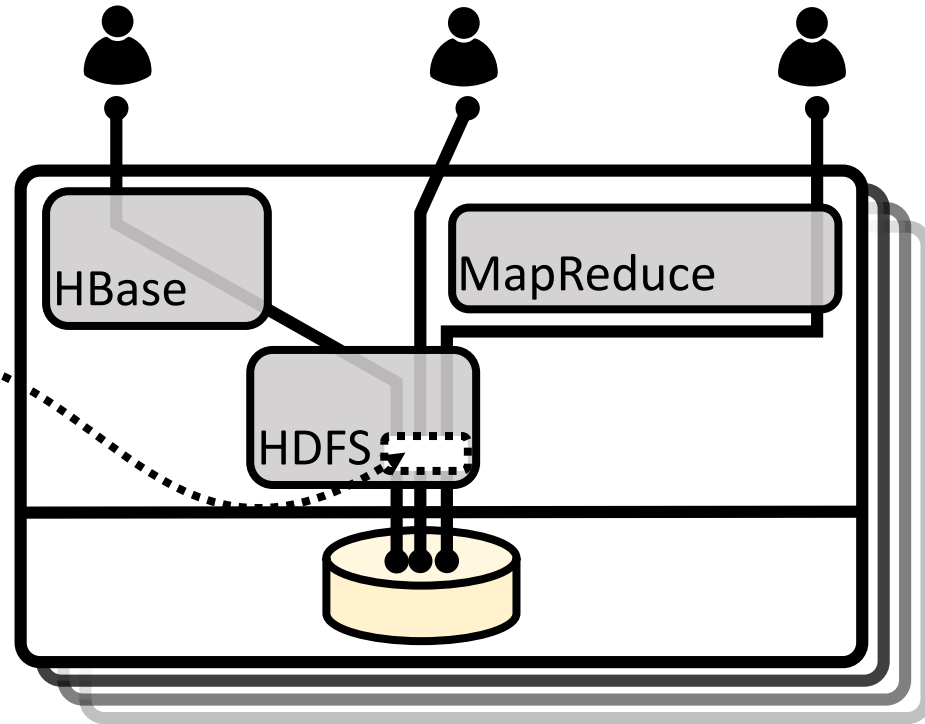
DataNode Metrics

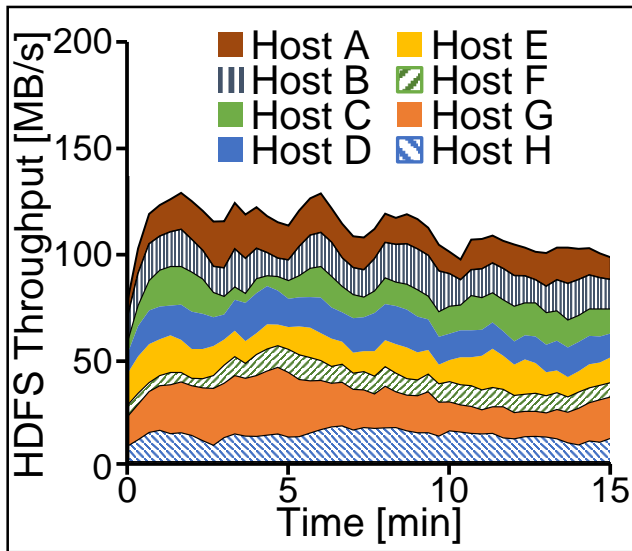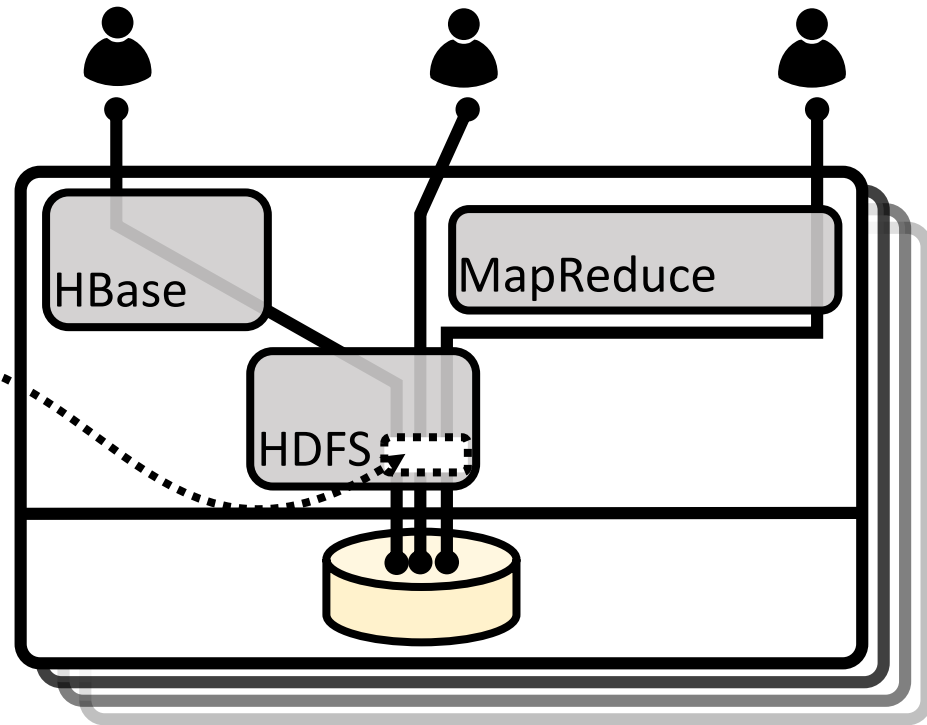("DataNodeMetrics", delta=10, host="hop01", ...)

HBase

MapReduce

HDFS

HDFS Throughput [MB/s]

- Host A
- Host B
- Host C
- Host D
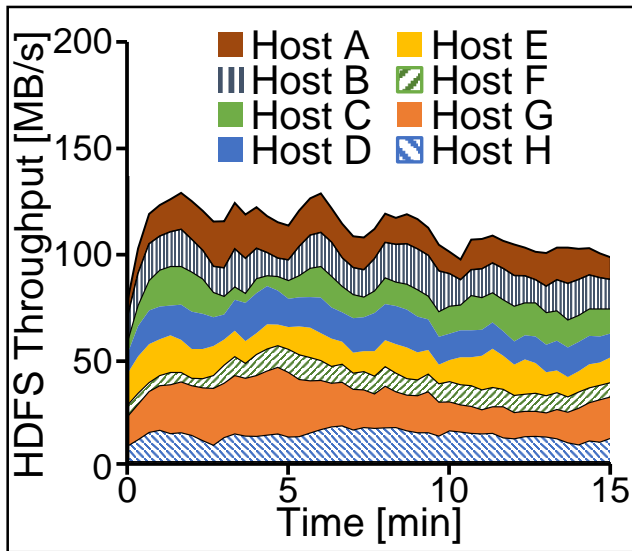- Host E
- Host F
- Host G
- Host H

Time [min]

```
From incr In DataNodeMetrics.incrBytesRead
  GroupBy incr.host
  Select incr.host, SUM(incr.delta)
```

15

DataNodeMetrics.java

```
50    public class DataNodeMetrics {
        ...
266     public void incrBytesRead(int delta) {
267         ...
268     }
      ...
407   }
```

("DataNodeMetrics", delta=10, host="hop01", …)

DataNode Metrics

HBase    MapReduce

HDFS

HDFS Throughput [MB/s]

■ Host A    ■ Host E
|||Host B    ▨ Host F
■ Host C    ■ Host G
■ Host D    ▨ Host H

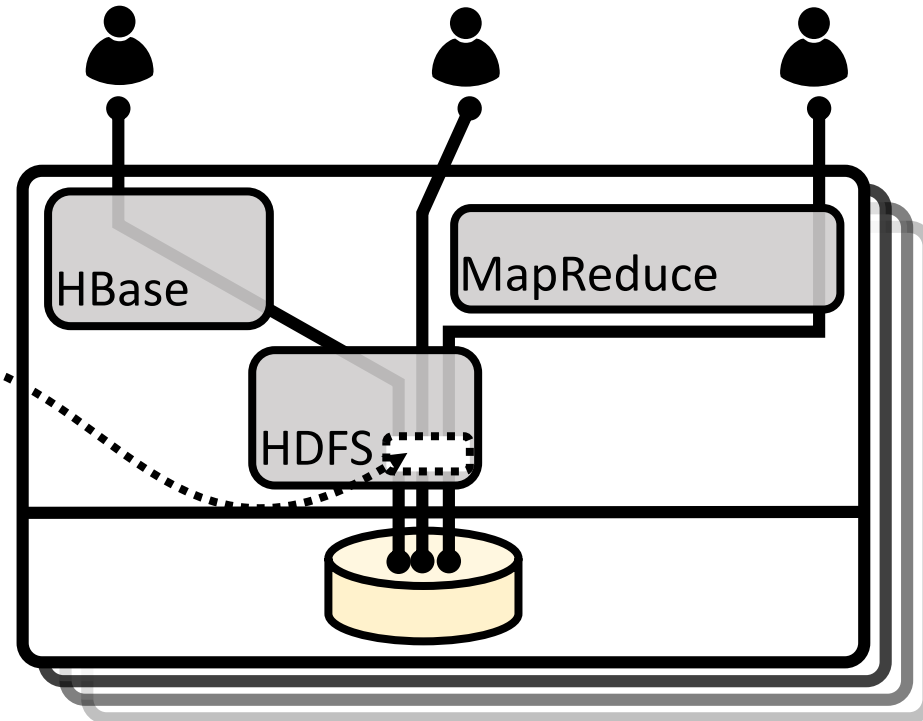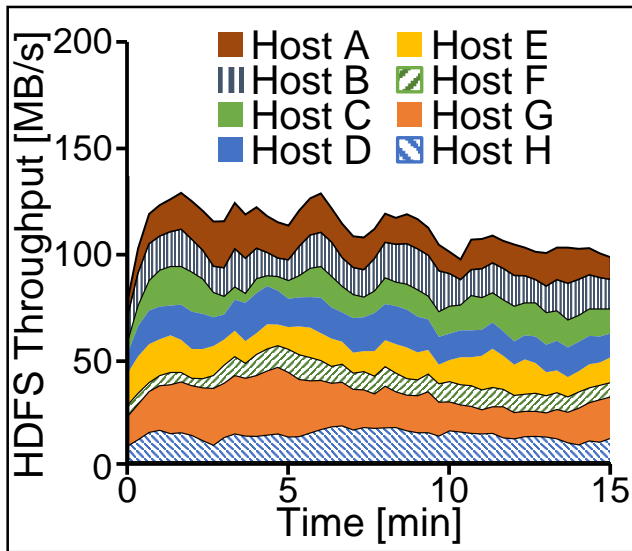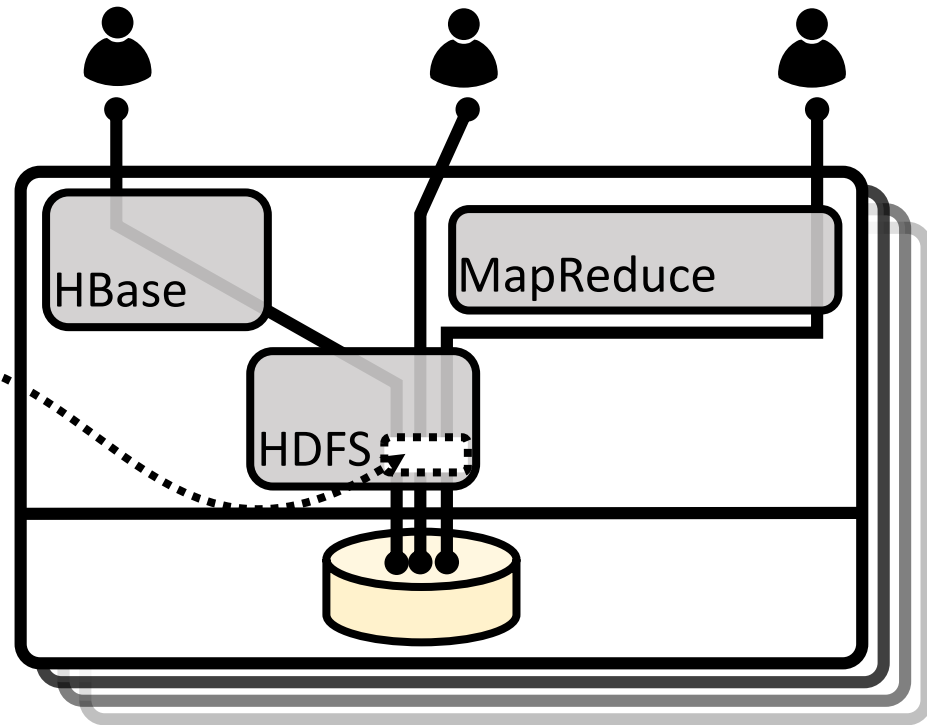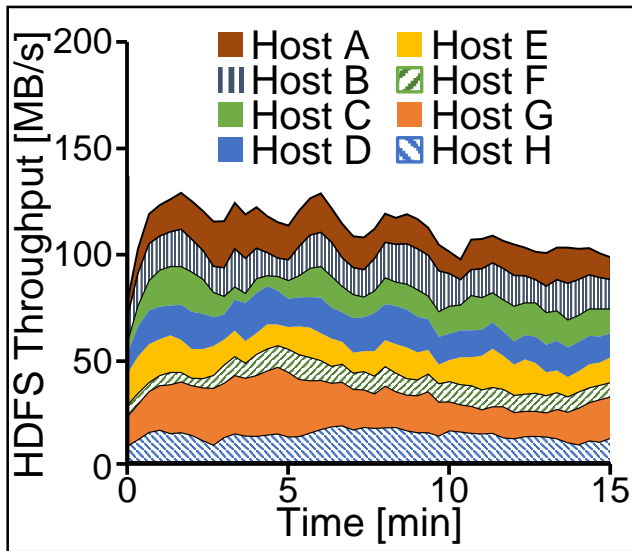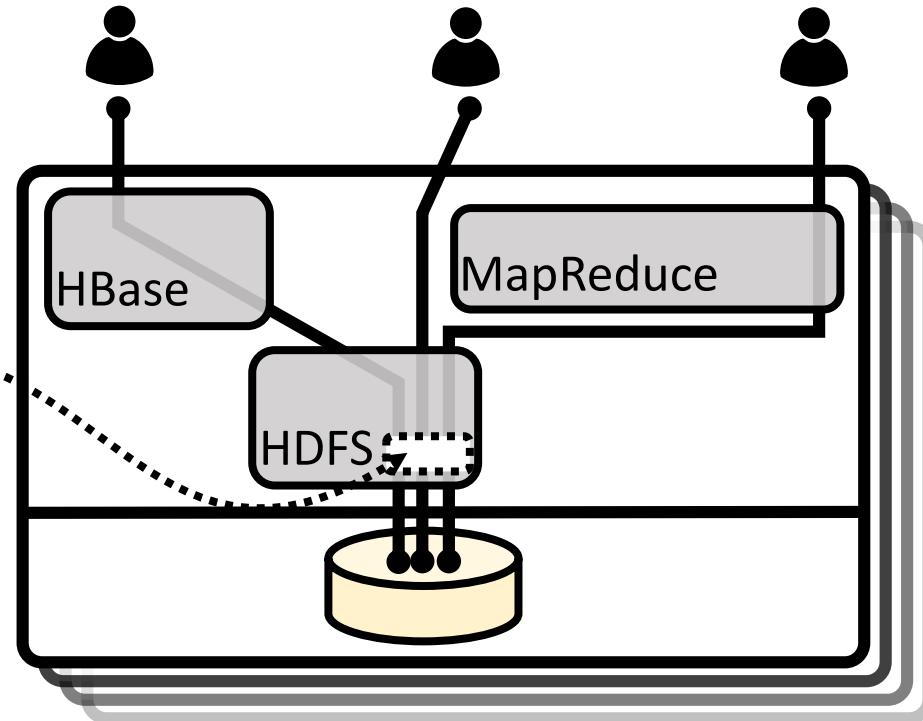200

150

100

50

0

0       5       10      15
Time [min]

From incr In DataNodeMetrics.incrBytesRead
  GroupBy incr.host
Select incr.host, SUM(incr.delta)

```
DataNodeMetrics.java
50    public class DataNodeMetrics {
        ...
266     public void incrBytesRead(int delta) {
267       ...
268     }
        ...
407   }
```

**Tracepoint**
  **Class:** DataNodeMetrics
  **Method:** incrBytesRead
  **Exports:** "delta"=delta

("DataNodeMetrics", delta=10, host="hop01", …)

```
From incr In DataNodeMetrics.incrBytesRead
  GroupBy incr.host
  Select incr.host, SUM(incr.delta)
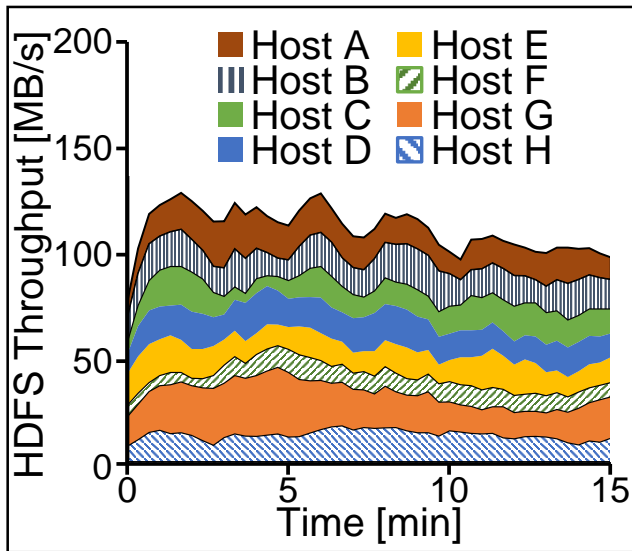```

Client Protocols

DataNode Metrics

HBase

MapReduce

HDFS



HDFS Throughput [MB/s]

- MRSORT10G
- MRSORT100G
- FSREAD64M
- FSREAD4M
- HGET
- HSCAN

Time [min]

# Happened-before Join (⋈→)

# Happened-before Join (⋈→)

# Happened-before Join ($\bowtie$)



("ClientProtocols", procName="HGET", …)

# Happened-before Join (⋈→)



Client Protocols

("ClientProtocols", procName="HGET", …)

DataNode Metrics

("DataNodeMetrics", delta=10, host="Hop01", …)

HBase

MapReduce

HDFS

HDFS Throughput [MB/s]

- MRSORT10G
- MRSORT100G
- FSREAD64M
- FSREAD4M
- HGET
- HSCAN

Time [min]

# Happened-before Join ($\bowtie$)



("ClientProtocols", procName="HGET", …)

("DataNodeMetrics", delta=10, host="Hop01", …)

# Happened-before Join (⋈)



("ClientProtocols", procName="HGET", …)

("DataNodeMetrics", delta=10, host="Hop01", …)

Client Protocols ⋈ DataNode Metrics

("ClientProtocols", procName="HGET", …
"DataNodeMetrics", delta=10, host="Hop01", …)

17

# Happened-before Join (⋈⃗)



("**ClientProtocols**", procName="HGET", …)

("DataNodeMetrics", delta=10, host="Hop01", …)

```
From incr In DataNodeMetrics.incrBytesRead
  Join client In First(ClientProtocols) On client -> incr
  GroupBy client.procName
  Select client.procName, SUM(incr.delta)
```
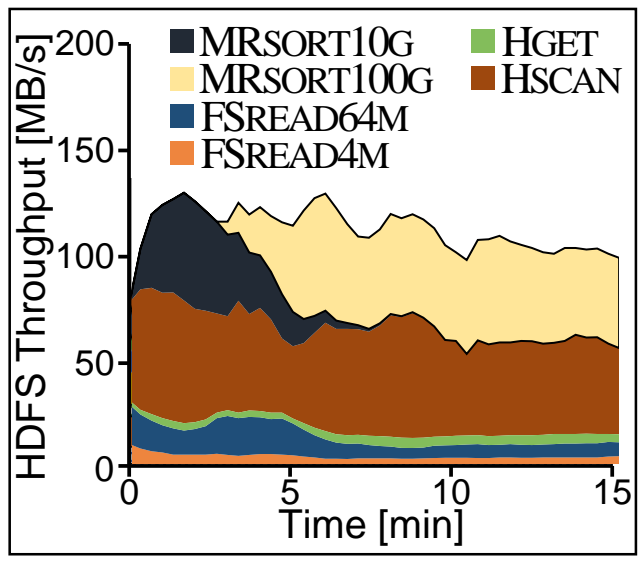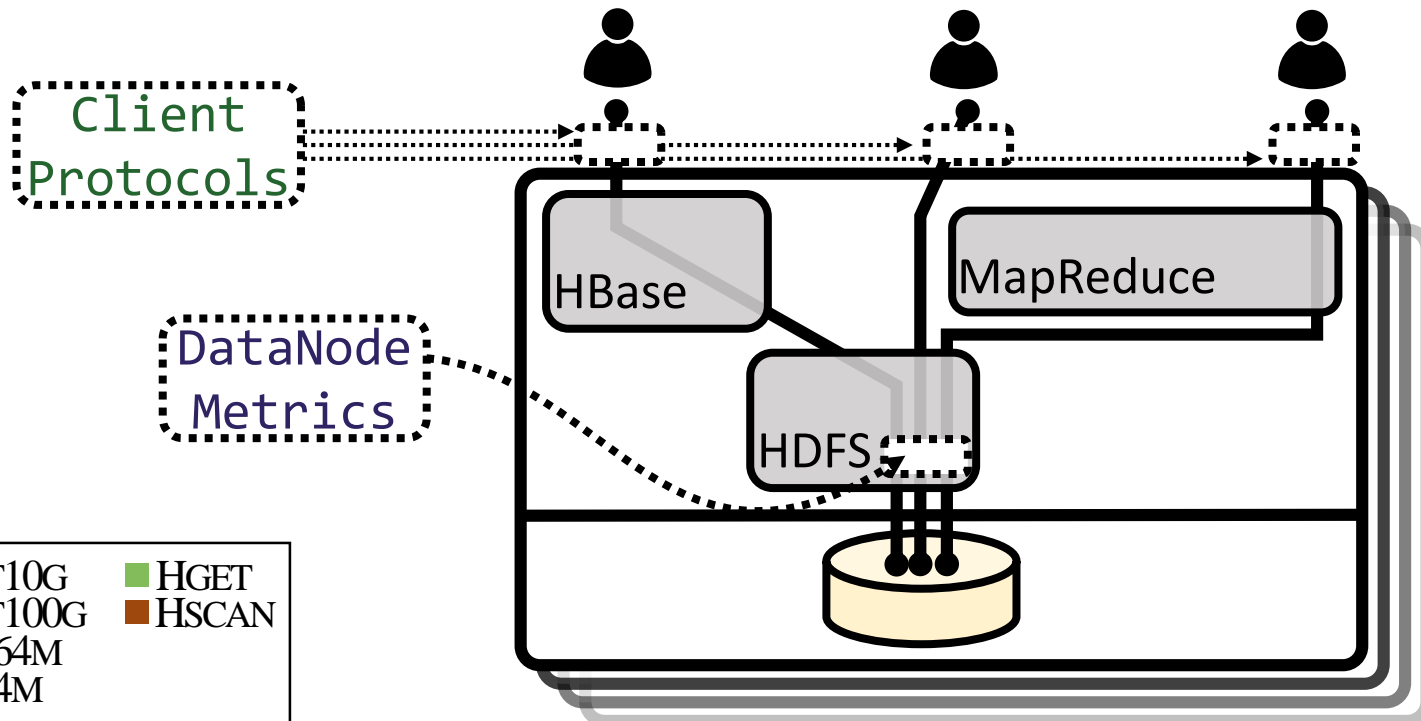
# Happened-before Join (⋈⃗)



Client
Protocols

("ClientProtocols", procName="HGET", …)

DataNode
Metrics

("DataNodeMetrics", delta=10, host="Hop01", …)

HBase

MapReduce

HDFS

```
From incr In DataNodeMetrics.incrBytesRead
  Join client In First(ClientProtocols) On client -> incr
  GroupBy client.procName
  Select client.procName, SUM(incr.delta)
```

# Happened-before Join (⋈⃗)



Client
Protocols

("**ClientProtocols**", procName="HGET", …)

DataNode
Metrics

("DataNodeMetrics", delta=10, host="Hop01", …)

HBase

MapReduce

HDFS

```
From incr In DataNodeMetrics.incrBytesRead
    Join client In First(ClientProtocols) On client -> incr
    GroupBy client.procName
    Select client.procName, SUM(incr.delta)
```

# Happened-before Join (⋈⃗)



Client
Protocols

("ClientProtocols", procName="HGET", …)

DataNode
Metrics

("DataNodeMetrics", delta=10, host="Hop01", …)

HBase

MapReduce

HDFS

```
From incr In DataNodeMetrics.incrBytesRead
  Join client In First(ClientProtocols) On client -> incr
  GroupBy client.procName
  Select client.procName, SUM(incr.delta)
```

# Happened-before Join (⋈⃗)



Client
Protocols

("ClientProtocols", procName="HGET", …)

DataNode
Metrics

("DataNodeMetrics", delta=10, host="Hop01", …)

HBase

MapReduce

HDFS

```
From incr In DataNodeMetrics.incrBytesRead
  Join client In First(ClientProtocols) On client -> incr
  GroupBy client.procName
  Select client.procName, SUM(incr.delta)
```

# Happened-before Join ( $\bowtie$ )



Client Protocols

("ClientProtocols", procName="HGET", ...)

DataNode Metrics

("DataNodeMetrics", delta=10, host="Hop01", ...)

HBase

MapReduce

HDFS

Legend:
- MRSORT10G
- MRSORT100G
- FSREAD64M
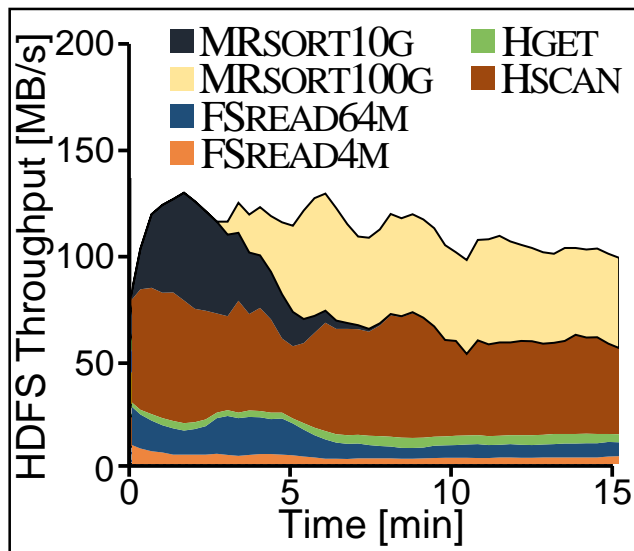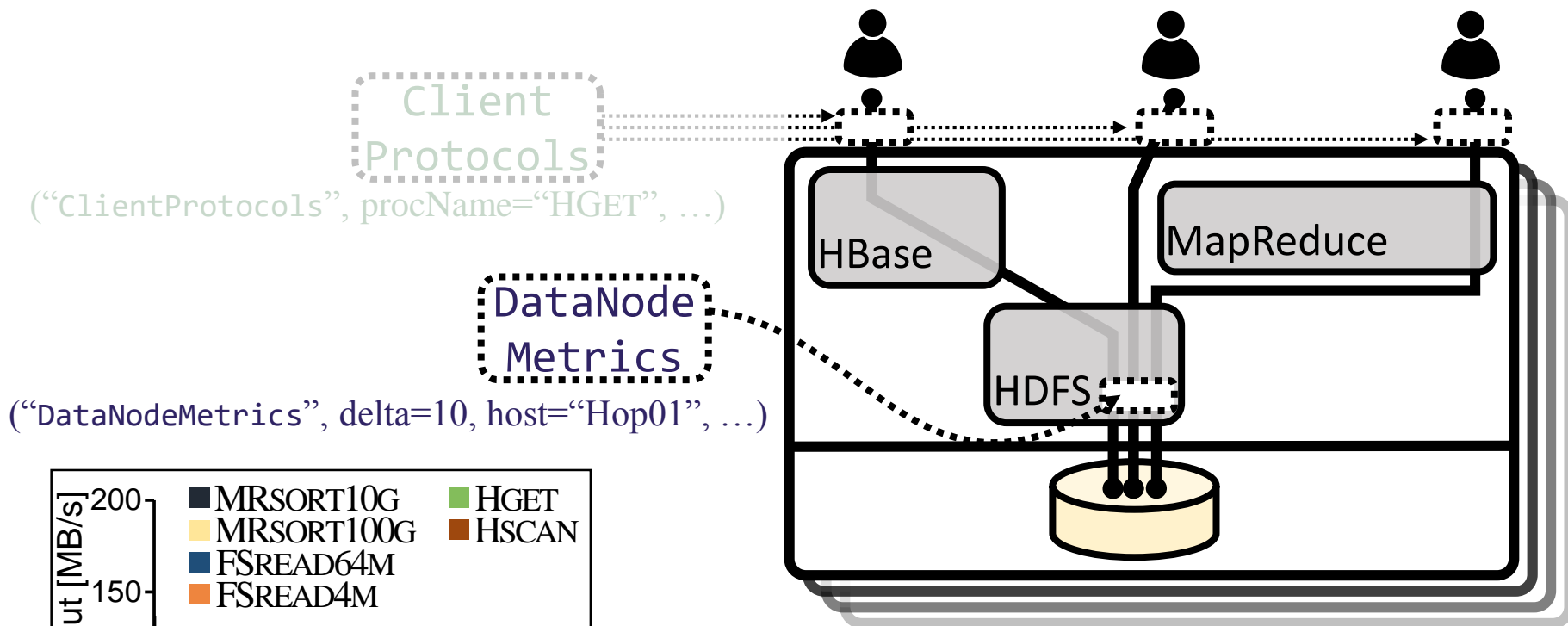- FSREAD4M
- HGET
- HSCAN

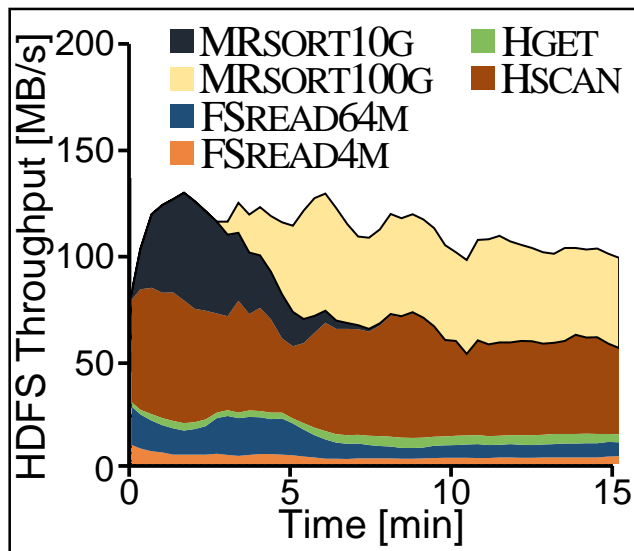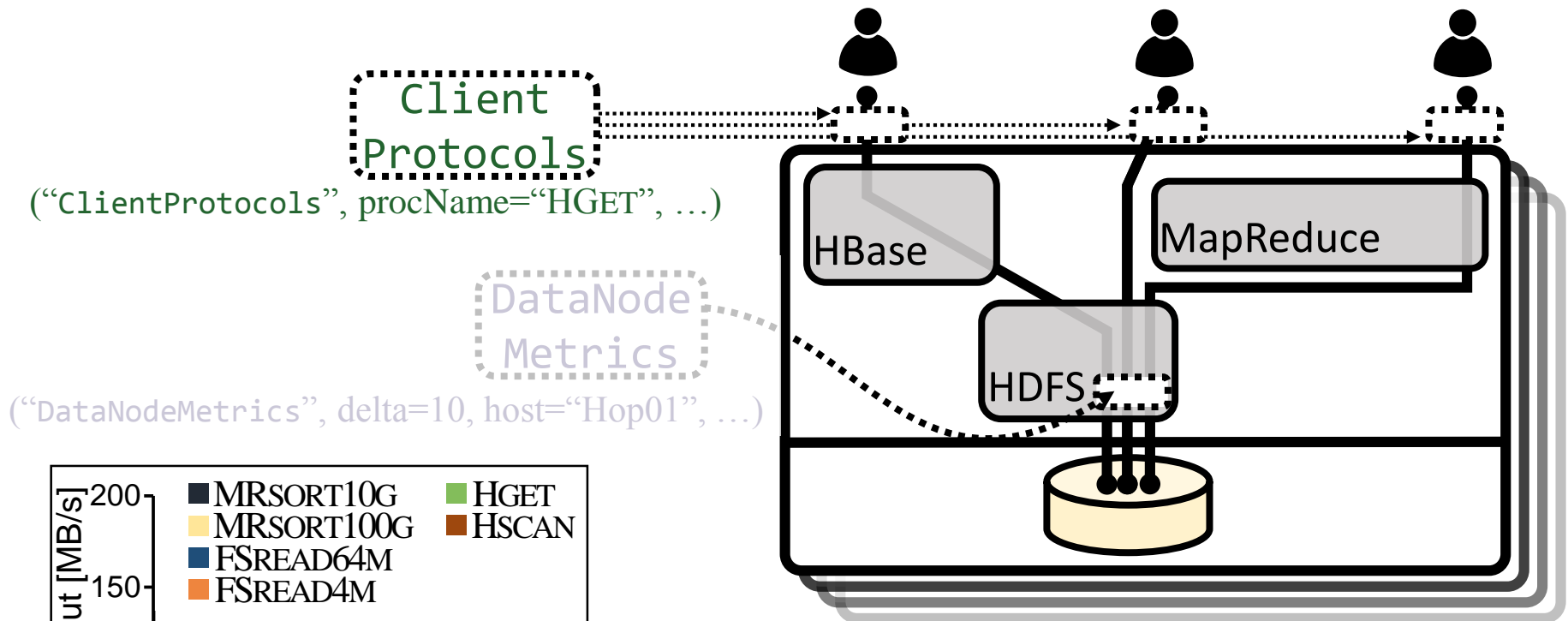HDFS Throughput [MB/s] vs Time [min]

```
From incr In DataNodeMetrics.incrBytesRead
  Join client In First(ClientProtocols) On client -> incr
  GroupBy client.procName
  Select client.procName, SUM(incr.delta)
```

(procName="HGET", delta=10, ...)

# Design & Implementation
## Pivot Tracing Pre-requisites

# Design & Implementation
## Pivot Tracing Pre-requisites

**Dynamic instrumentation**          PT Agent

# Design & Implementation
## Pivot Tracing Pre-requisites

**Dynamic instrumentation** ⟿ PT Agent

**Causal tracing** ⟿ Baggage

**Causal tracing** $\rightsquigarrow$ Baggage

Legend:
- Process
- Tenant
- Workload
- Request

Baggage is a Key:Value container propagated alongside a request

- Generalization of metadata in end-to-end tracing
- One instance per request

**Causal tracing** ⟿ Baggage

Baggage is a Key:Value container propagated alongside a request

- Generalization of metadata in end-to-end tracing
- One instance per request

**Causal tracing** ⟿ 🧳 Baggage

Legend:
- Process
- Tenant
- Workload
- Request

Baggage is a Key:Value container propagated alongside a request

- Generalization of metadata in end-to-end tracing
- One instance per request

**Causal tracing** ⟿ Baggage

Baggage is a Key:Value container propagated alongside a request

- Generalization of metadata in end-to-end tracing
- One instance per request

**Causal tracing** ⟿ Baggage

Legend:
- Process
- Tenant
- Workload
- Request

clientName="FSRead"

🧳 Baggage is a Key:Value container propagated alongside a request

- Generalization of metadata in end-to-end tracing
- One instance per request

**Causal tracing** ⟿ 🧳 Baggage

Legend:
- Process
- Tenant
- Workload
- Request

clientName="FSRead"

UNPACK(clientName)

clientName="FSRead"

🧳 Baggage is a Key:Value container propagated alongside a request

- Generalization of metadata in end-to-end tracing
- One instance per request

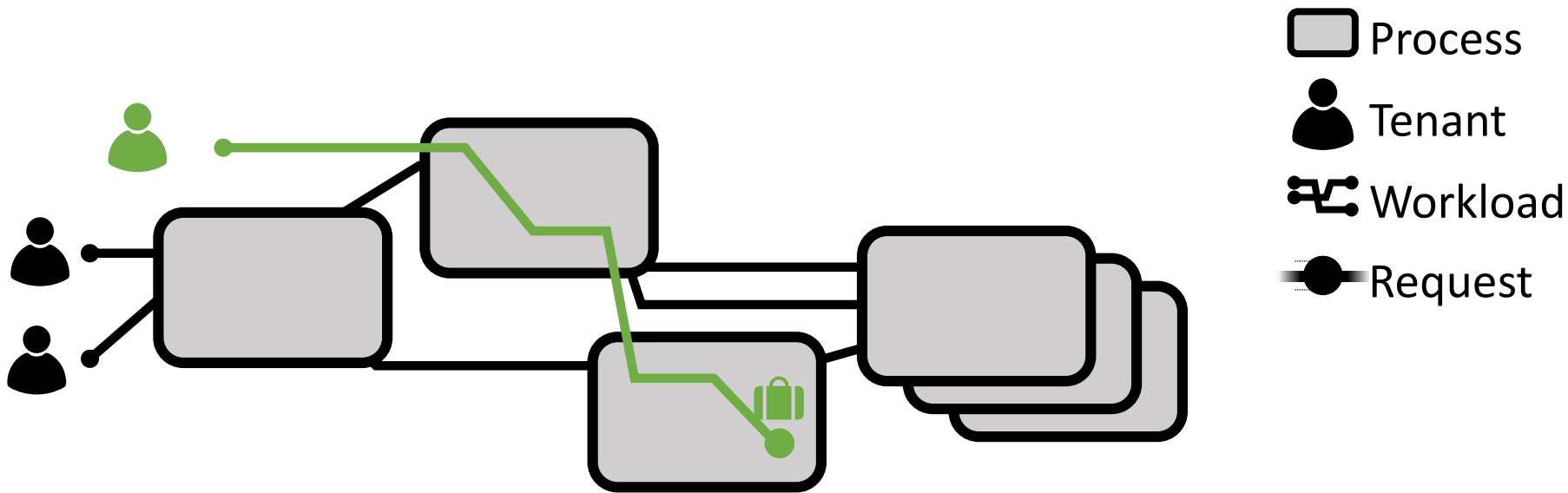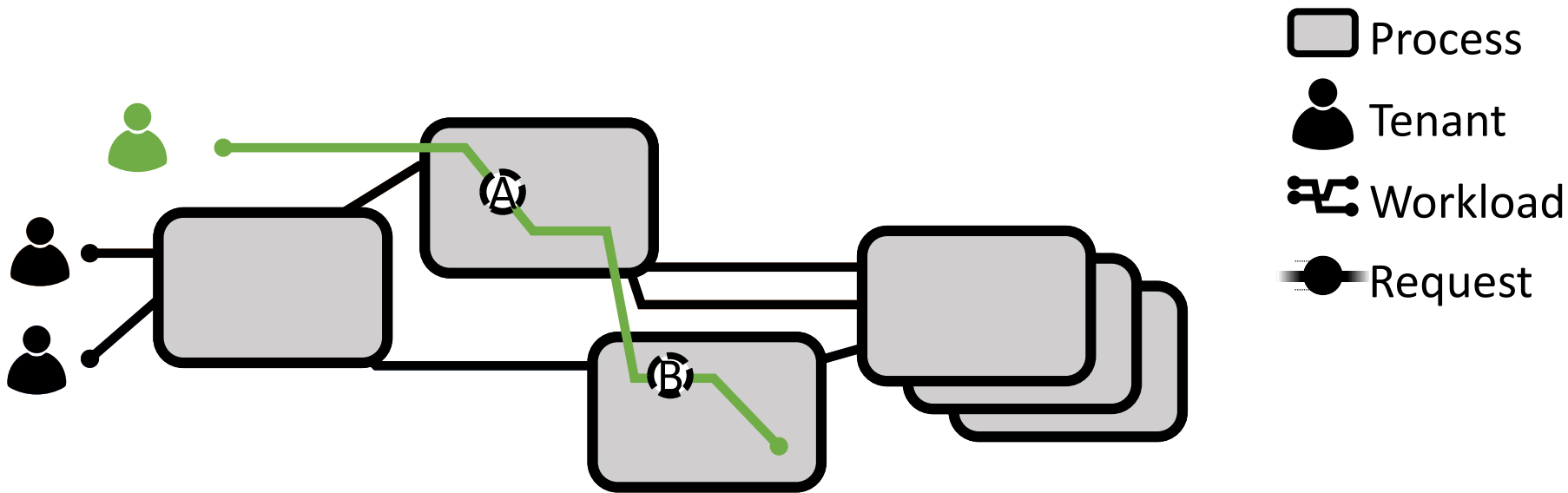**Causal tracing** ⟿ 🧳 Baggage

21

Process
Tenant
Workload
Request

21

Process
Tenant
Workload
Request

+ 🧳 Baggage
+ 🔧 PT Agent

Pivot Tracing Enabled

# Design & Implementation
## Queries

PT Agent

PT Agent

# Tracepoints

Places where PT can add instrumentation

# Tracepoints

Places where PT can add instrumentation

**Tracepoint A**
 **Class:** A
 **Method:** A1()

**Tracepoint B**
 **Class:** B
 **Method**: B1()
 **Exports:** "delta"=delta

# Tracepoints

Places where PT can add instrumentation

Export identifiers accessible to queries
Defaults: host, timestamp, pid, proc name

**Tracepoint A**
 **Class:** A
 **Method:** A1()

**Tracepoint B**
 **Class:** B
 **Method:** B1()
 **Exports:** "delta"=delta

# Tracepoints

Places where PT can add instrumentation

Export identifiers accessible to queries
 Defaults: host, timestamp, pid, proc name

Only references – not materialized until query is installed

**Tracepoint A**
 **Class:** A
 **Method:**  A1()

**Tracepoint B**
 **Class:** B
 **Method**: B1()
 **Exports:** "delta"=delta

# Query Language

Relational query language, similar to SQL,
LINQ

- Selection
- Projection
- Filter

- GroupBy
- Aggregation
- Happened-Before Join

Refers to tracepoint-exported identifiers

**Tracepoint A**
 **Class:** A
 **Method:** A1()

**Tracepoint B**
 **Class:** B
 **Method**: B1()
 **Exports:** "delta"=delta

26

# Query Language

Relational query language, similar to SQL, LINQ

- Selection
- Projection
- Filter

- GroupBy
- Aggregation
- Happened-Before Join

Refers to tracepoint-exported identifiers

```
From a In A
    Join b In B On a -> b
    GroupBy a.procName
    Select a.procName, SUM(b.delta)
```

**Tracepoint A**
  **Class:** A
  **Method:** A1()

**Tracepoint B**
  **Class:** B
  **Method:** B1()
  **Exports:** "delta"=delta

# Query Language

Relational query language, similar to SQL, LINQ

- Selection
- Projection
- Filter
- GroupBy
- Aggregation
- Happened-Before Join

Refers to tracepoint-exported identifiers

Output: stream of tuples

e.g., (procName, delta)

**PT**

```
From a In A
  Join b In B On a -> b
  GroupBy a.procName
  Select a.procName, SUM(b.delta)
```

**Tracepoint A**
 **Class:** A
 **Method:** A1()

**Tracepoint B**
 **Class:** B
 **Method:** B1()
 **Exports:** "delta"=delta

26

# Advice

Query is compiled to advice
(intermediate representation for instrumentation)

Pivot Tracing
Front End

Query

PT

```
From a In A
  Join b In B On a -> b
  GroupBy a.procName
  Select a.procName, SUM(b.delta)
```

**Advice A1**
OBSERVE    procName
PACK       procName

**Advice B1**
OBSERVE    delta
UNPACK     procName
EMIT       procName, SUM(delta)

27

# Advice

Query is compiled to advice
(intermediate representation for instrumentation)

Advice will be installed at tracepoints

Pivot Tracing
Front End

Query

PT

```
From a In A
  Join b In B On a -> b
  GroupBy a.procName
  Select a.procName, SUM(b.delta)
```

**Advice A1**                                     A
OBSERVE   procName
PACK      procName

**Advice B1**                                     B
OBSERVE   delta
UNPACK    procName
EMIT      procName, SUM(delta)

27

## Advice

Query is compiled to advice
(intermediate representation for instrumentation)

Advice will be installed at tracepoints

Limited instruction set
   OBSERVE
   PACK
   FILTER
   UNPACK
   EMIT

**Pivot Tracing Front End**

Query

**PT**

```
From a In A
   Join b In B On a -> b
   GroupBy a.procName
   Select a.procName, SUM(b.delta)
```

**Advice A1**                                    Ⓐ
OBSERVE   procName
PACK      procName

**Advice B1**                                    Ⓑ
OBSERVE   delta
UNPACK    procName
EMIT      procName, SUM(delta)

27

Weaving

PT Agent dynamically enables
advice at tracepoints

```
From a In A
    Join b In B On a -> b
    GroupBy a.procName
    Select a.procName, SUM(b.delta)
```

**Advice A1**  A
```
OBSERVE   procName
PACK      procName
```

**Advice B1**  B
```
OBSERVE   delta
UNPACK    procName
EMIT      procName, SUM(delta)
```

28

Instrumented System   (+ 🧳 Baggage, + 🔧 PT Agent)

A1

B1

PT Agent

PT Agent

Pivot Tracing
Front End

Advice A1

Advice B1

Query

PT

Evaluating ⋈

```
From a In A
    Join b In B On a -> b
    GroupBy a.procName
    Select a.procName, SUM(b.delta)
```

**Advice A1**                                    A
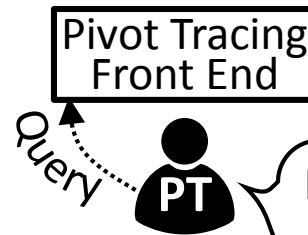```
OBSERVE   procName
PACK      procName
```

**Advice B1**                                    B
```
OBSERVE   delta
UNPACK    procName
EMIT      procName, SUM(delta)
```

29

Instrumented System   (+ 🧳 Baggage, + 🔧 PT Agent)

Evaluating ⋈

PT Agent

B1

PT Agent

Pivot Tracing
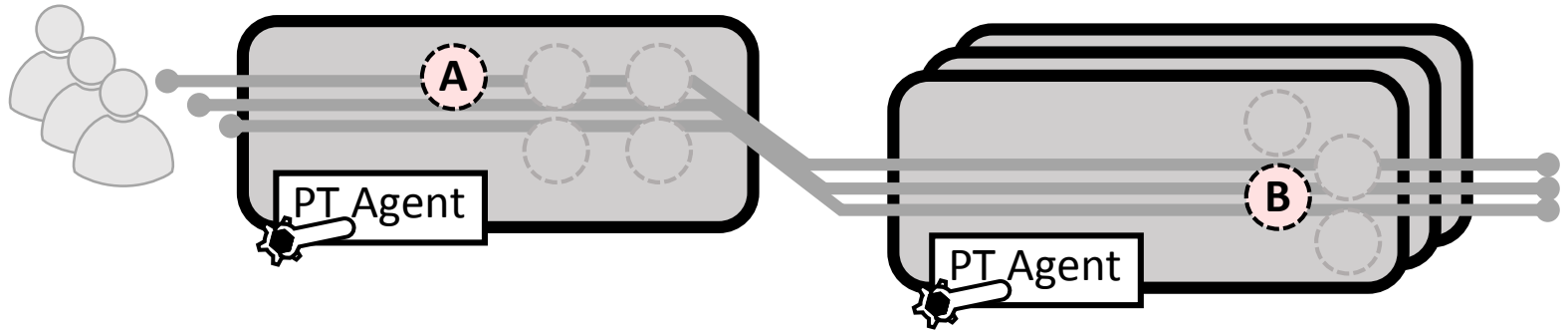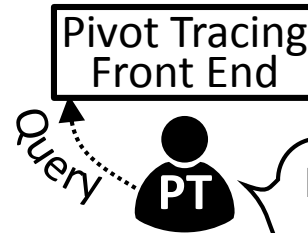Front End

Advice A1

Advice B1

Query

PT

```
From a In A
    Join b In B On a -> b
GroupBy a.procName
Select a.procName, SUM(b.delta)
```

**Advice A1**    A
OBSERVE    procName
PACK       procName

**Advice B1**    B
OBSERVE  delta
UNPACK   procName
EMIT     procName, SUM(delta)

Instrumented System   (+ 🧳 Baggage, + 🔧 PT Agent)

OBSERVE procName

PT Agent

B1

PT Agent

Pivot Tracing
Front End

Advice A1

Advice B1

Query
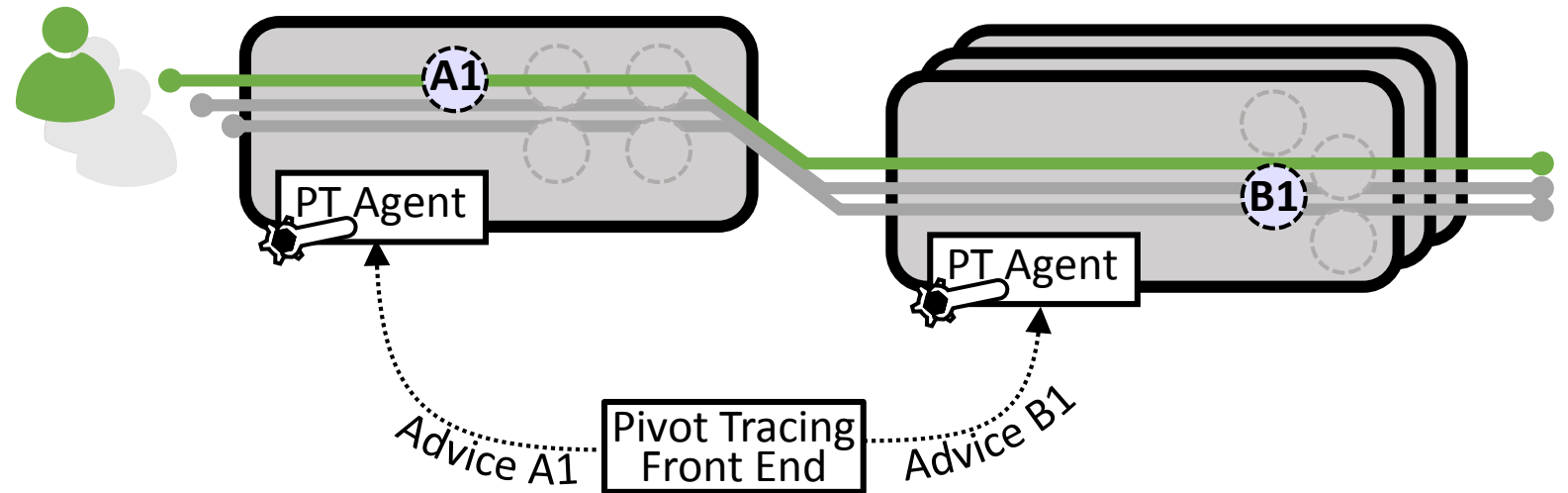
PT

Evaluating ⋈

```
From a In A
    Join b In B On a -> b
    GroupBy a.procName
    Select a.procName, SUM(b.delta)
```

**Advice A1**                    **A**
OBSERVE   procName
PACK      procName

**Advice B1**                    B
OBSERVE  delta
UNPACK   procName
EMIT     procName, SUM(delta)

procName

A1 PACK

PT Agent

B1

PT Agent

Advice A1

Pivot Tracing
Front End

Advice B1

Query

PT

```
From a In A
  Join b In B On a -> b
  GroupBy a.procName
  Select a.procName, SUM(b.delta)
```

Evaluating ⋈

**Advice A1**                                     A
OBSERVE   procName
PACK      procName

**Advice B1**                                     B
OBSERVE  delta
UNPACK   procName
EMIT     procName, SUM(delta)

Instrumented System    (+ 🧳 Baggage, + 🔧 PT Agent)

A1

procName

PT Agent

B1

PT Agent

Pivot Tracing
Front End

Advice A1

Advice B1

Query

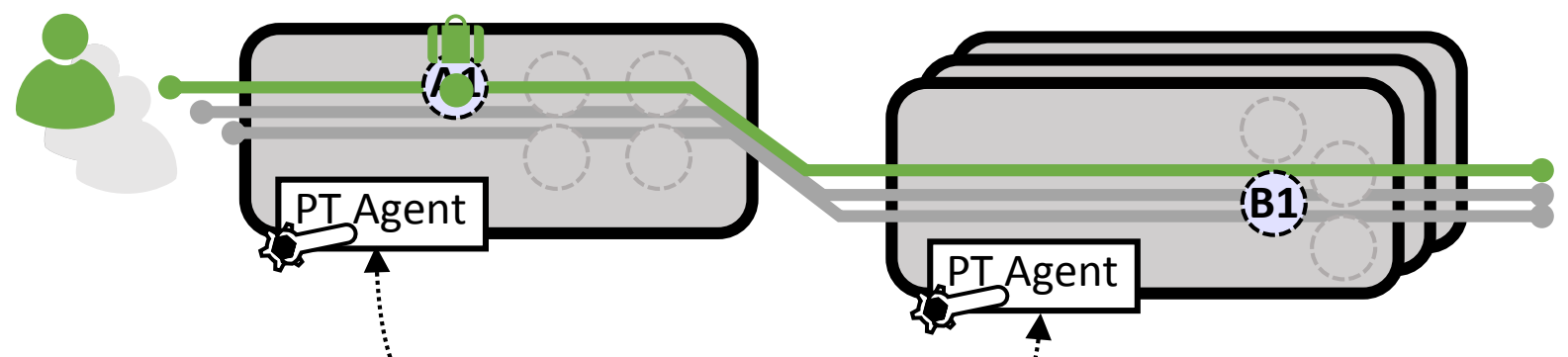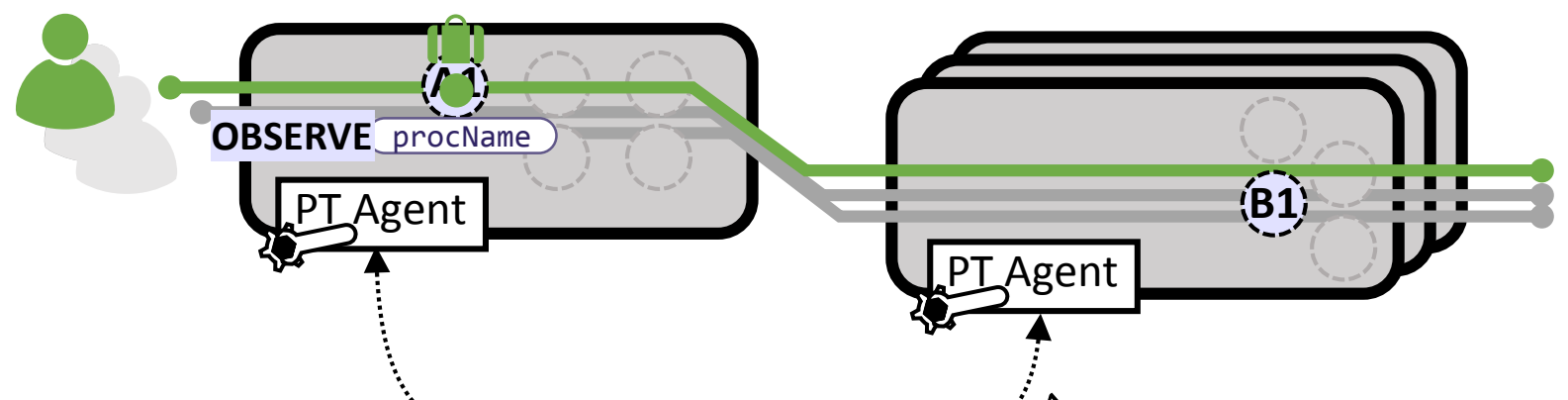Evaluating ⋈

PT

```
From a In A
  Join b In B On a -> b
  GroupBy a.procName
  Select a.procName, SUM(b.delta)
```

Advice A1                                              A
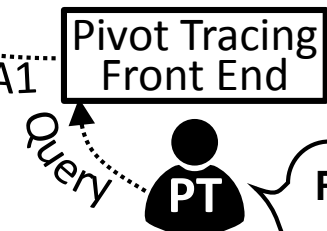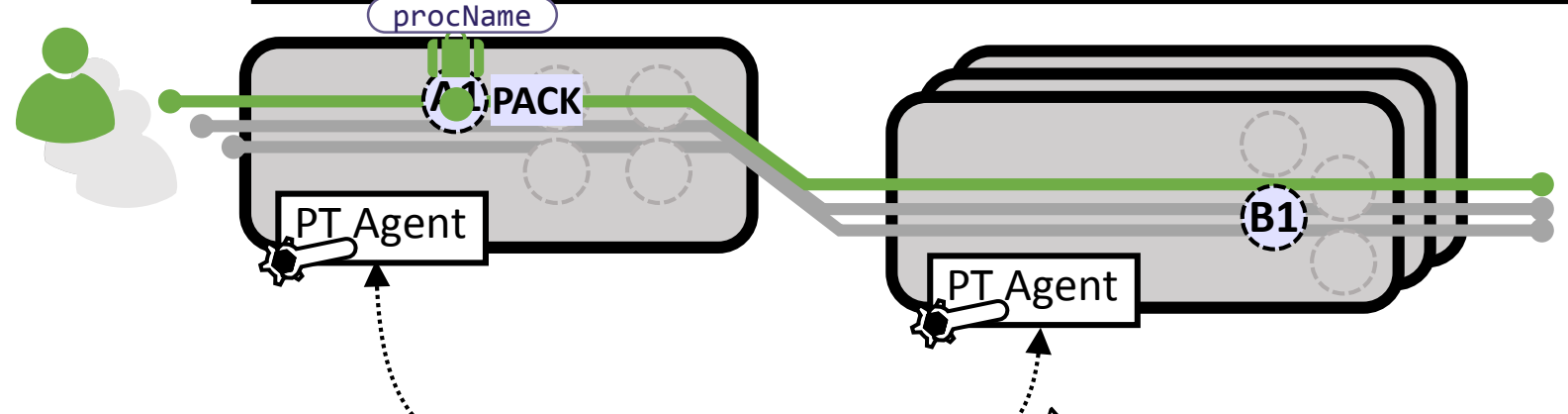OBSERVE    procName
PACK       procName

Advice B1                                              B
OBSERVE    delta
UNPACK     procName
EMIT       procName, SUM(delta)

Instrumented System    (+ 🧳 Baggage, + 🔧 PT Agent)

Evaluating ⋈

A1

procName

PT Agent

OBSERVE

B1

delta

PT Agent

Advice A1

Pivot Tracing
Front End

Advice B1

Query

PT

```
From a In A
    Join b In B On a -> b
    GroupBy a.procName
    Select a.procName, SUM(b.delta)
```

Advice A1                          A
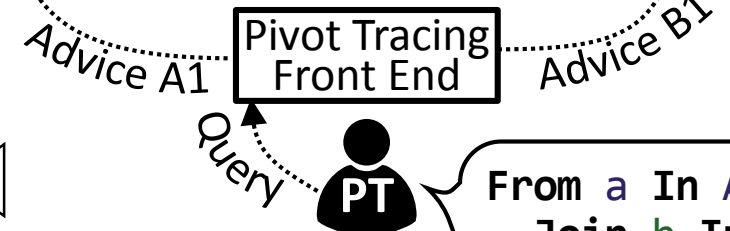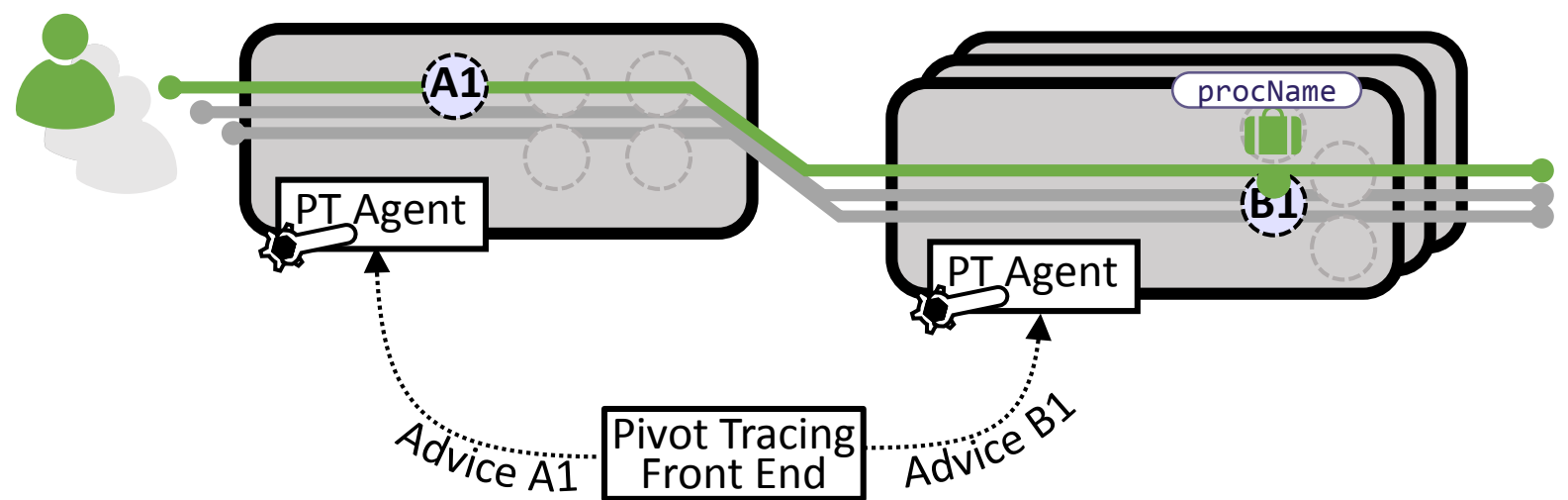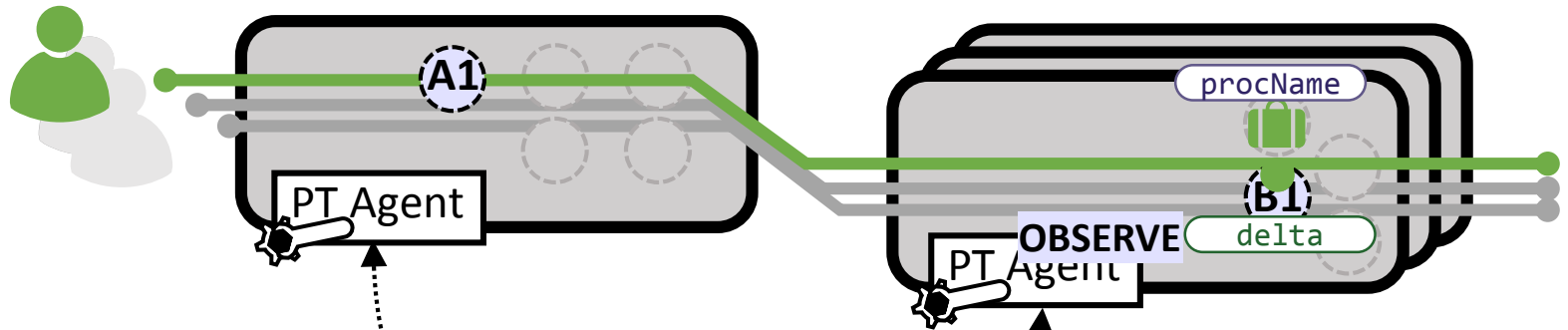OBSERVE    procName
PACK       procName

Advice B1                          B
OBSERVE    delta
UNPACK     procName
EMIT       procName, SUM(delta)

Instrumented System (+ 🧳 Baggage, + 🔧 PT Agent)

A1

procName

UNPACK

delta

procName

PT Agent

PT Agent

Pivot Tracing
Front End

Advice A1

Advice B1

Query

Evaluating ⋈

PT

```
From a In A
  Join b In B On a -> b
  GroupBy a.procName
  Select a.procName, SUM(b.delta)
```
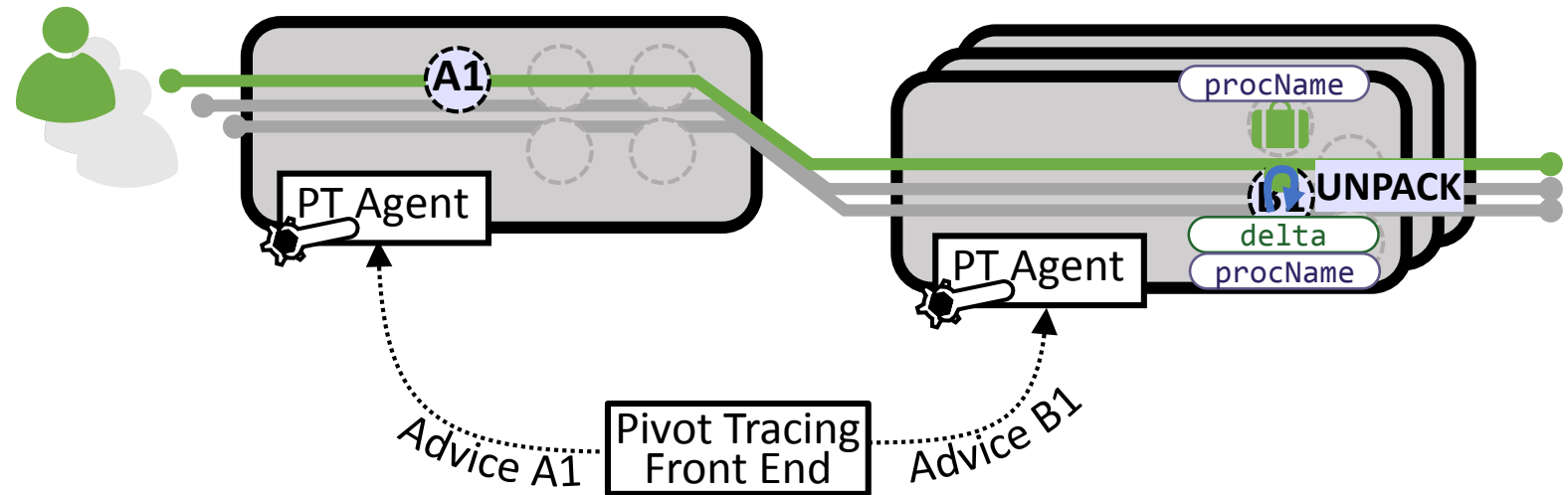
Advice A1                                    A
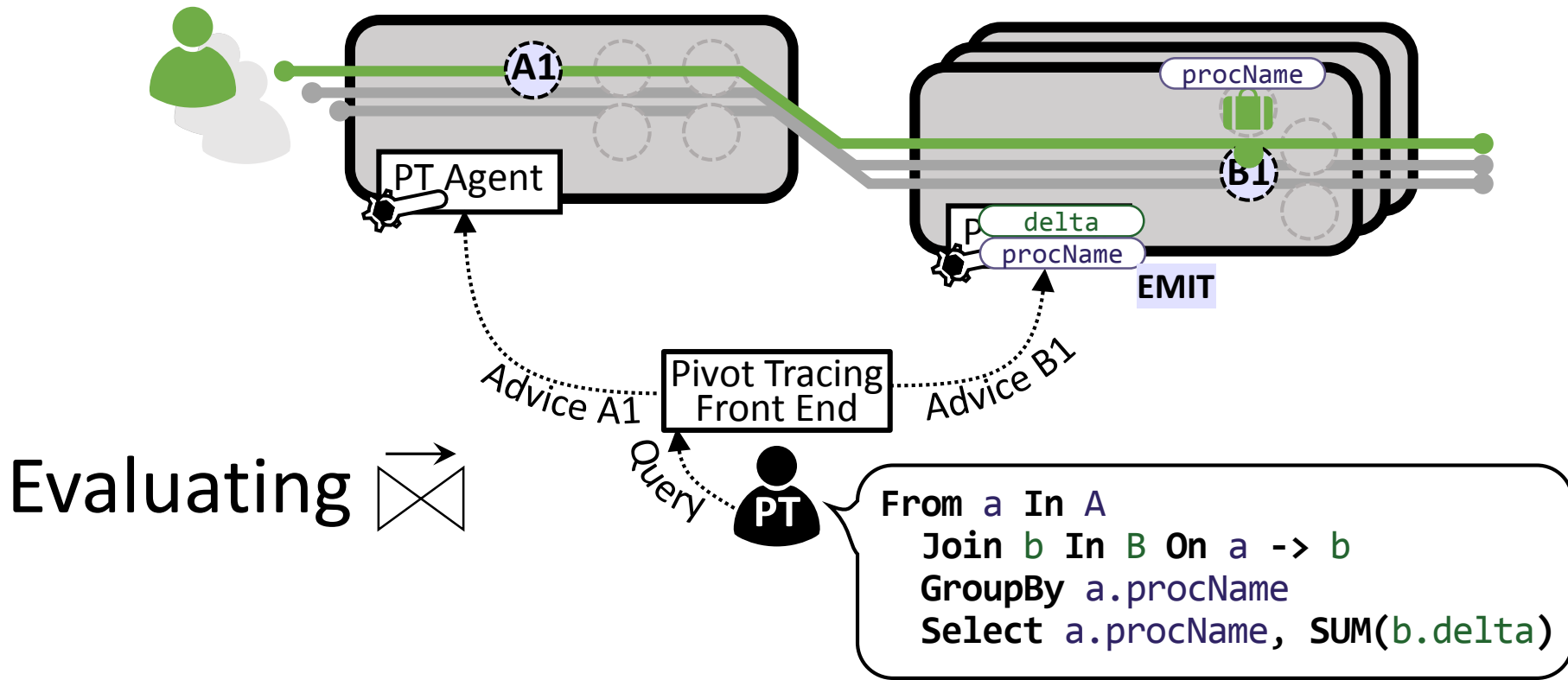OBSERVE   procName
PACK      procName

Advice B1                                    B
OBSERVE   delta
UNPACK    procName
EMIT      procName, SUM(delta)

Instrumented System    (+ 🧳 Baggage, + 🔧 PT Agent)

A1

procName

PT Agent

B1

delta
procName

P

EMIT

Pivot Tracing
Front End

Advice A1

Advice B1

Query

Evaluating ⋈

PT

```
From a In A
  Join b In B On a -> b
  GroupBy a.procName
  Select a.procName, SUM(b.delta)
```
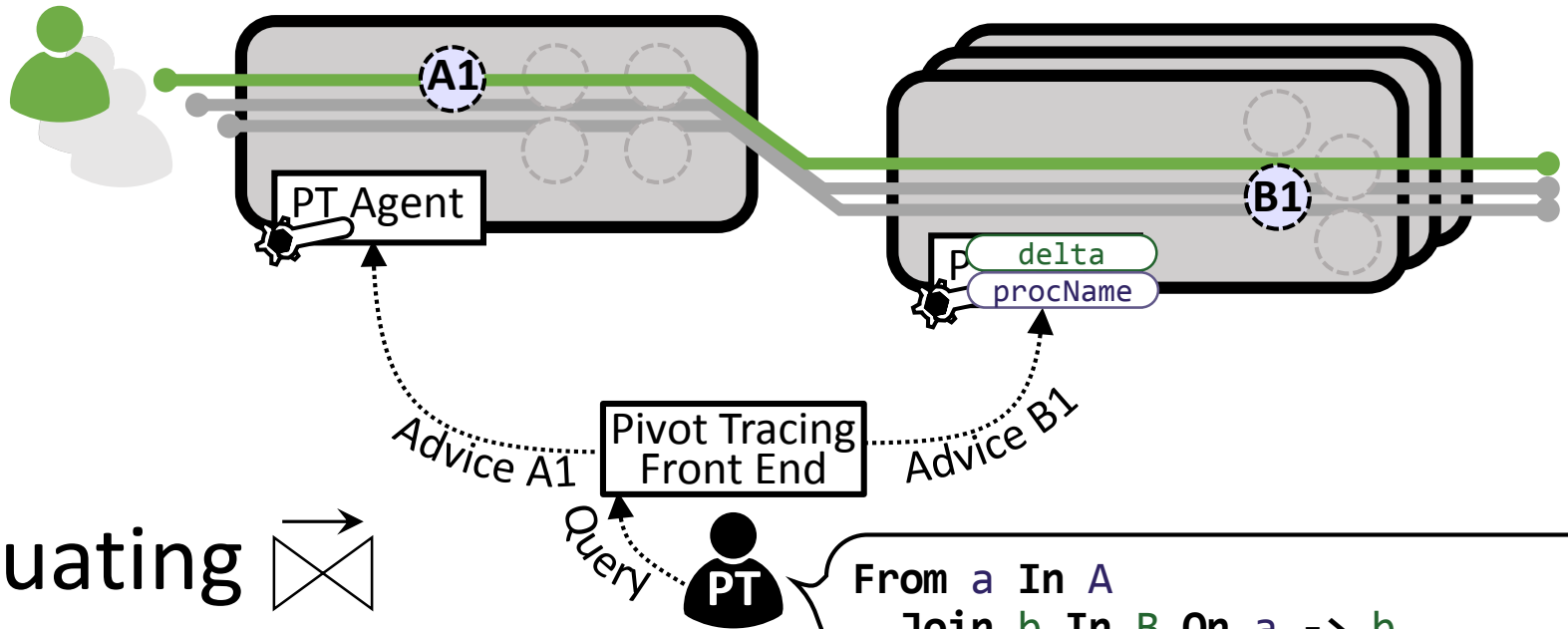
Advice A1                                    A
OBSERVE    procName
PACK       procName

Advice B1                                    B
OBSERVE    delta
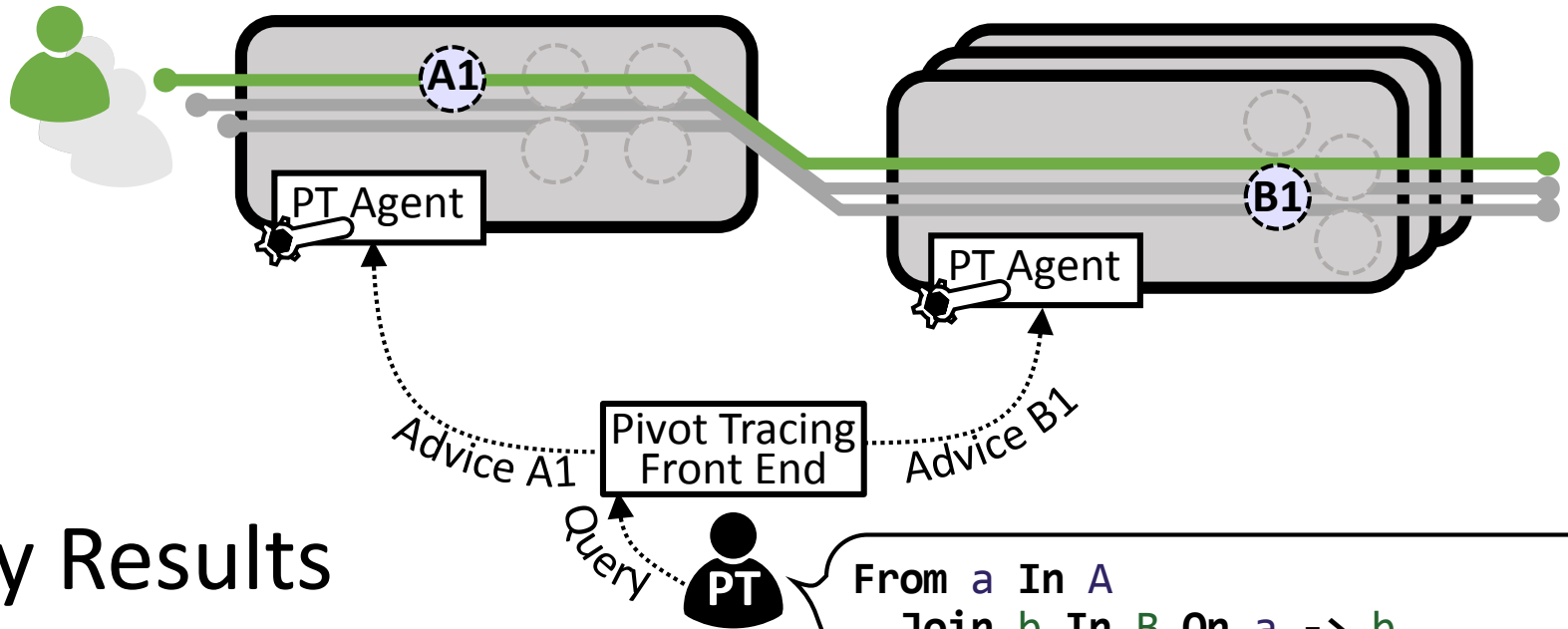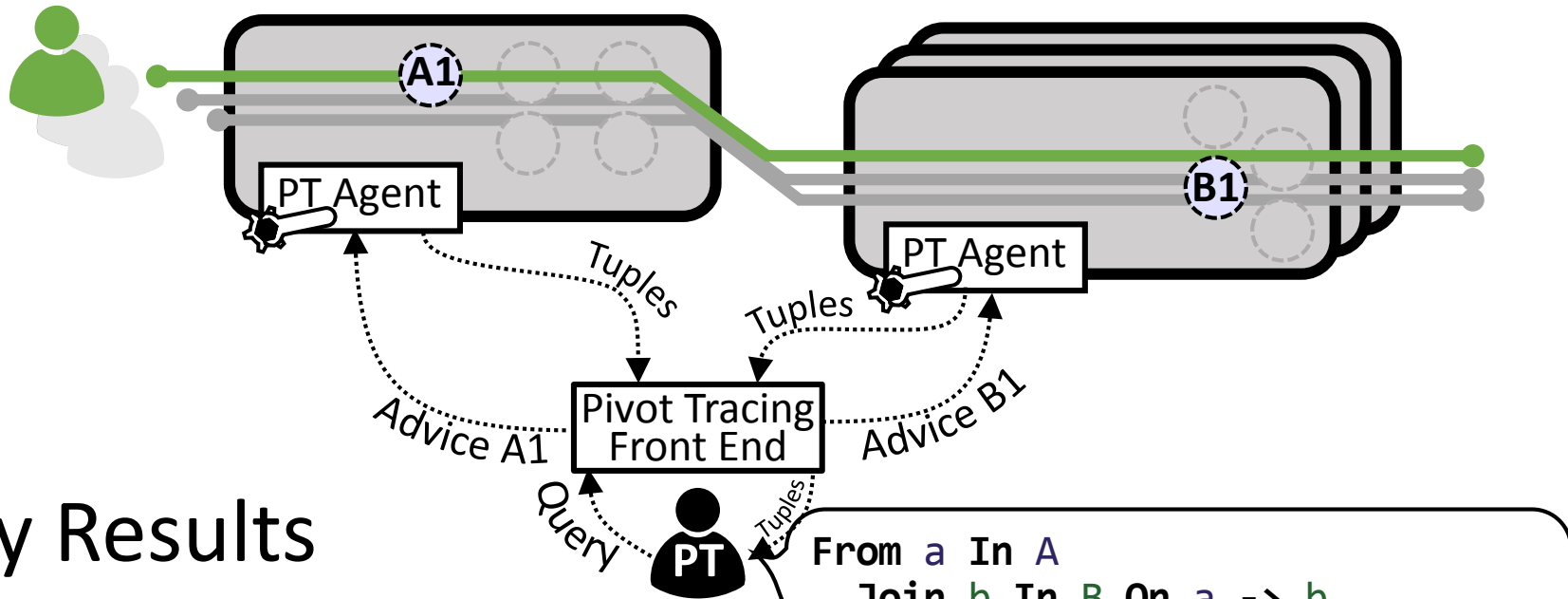UNPACK     procName
EMIT       procName, SUM(delta)

29

Instrumented System    (+ 🧳 Baggage, + 🔧 PT Agent)

A1

B1

PT Agent

delta
procName

Advice A1

Pivot Tracing
Front End

Advice B1

Query

PT

```
From a In A
    Join b In B On a -> b
    GroupBy a.procName
    Select a.procName, SUM(b.delta)
```

# Evaluating ⋈

Baggage explicitly follows execution

Evaluated inline during a request
 (no global aggregation needed)

**Advice A1**    A
OBSERVE    procName
PACK       procName

**Advice B1**    B
OBSERVE    delta
UNPACK     procName
EMIT       procName, SUM(delta)

29

# Query Results

Tuples are accumulated locally in PT Agent

Periodically reported back to user
 e.g., every second

PT Agent

A1

B1

PT Agent

Tuples

Tuples

Advice A1

Advice B1

Pivot Tracing
Front End

Query

Tuples

PT

```
From a In A
    Join b In B On a -> b
    GroupBy a.procName
    Select a.procName, SUM(b.delta)
```

**Advice A1**                    A
OBSERVE   procName
PACK      procName

**Advice B1**                    B
OBSERVE   delta
UNPACK    procName
EMIT      procName, SUM(delta)

29

# Pivot Tracing
## Evaluation

# Java-Based Implementation

# Java-Based Implementation

| PT Agent | PT agent thread that runs inside each process

- Javassist for dynamic instrumentation
- PubSub to receive commands / send tuples
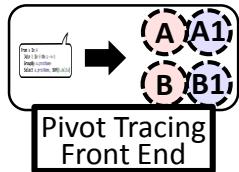
# Java-Based Implementation

PT Agent — PT agent thread that runs inside each process
- Javassist for dynamic instrumentation
- PubSub to receive commands / send tuples
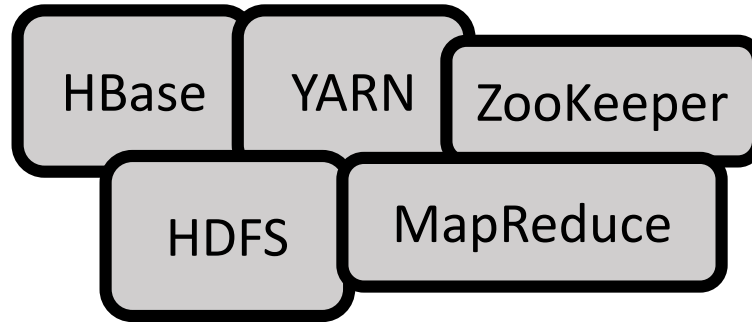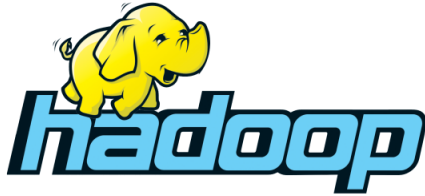
Baggage library for use by instrumented system
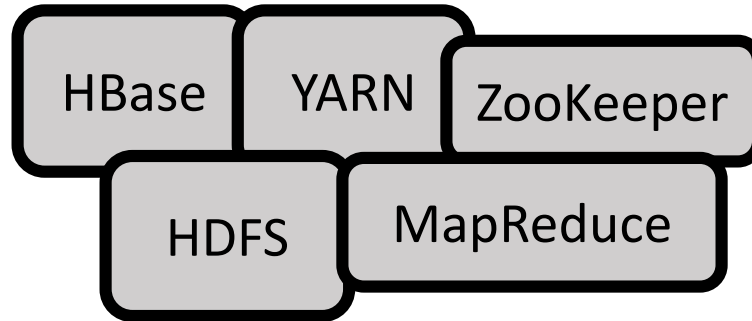- Data format specified using Protocol Buffers

# Java-Based Implementation

 PT agent thread that runs inside each process
- Javassist for dynamic instrumentation
- PubSub to receive commands / send tuples

 Baggage library for use by instrumented system
- Data format specified using Protocol Buffers

 Front-end client library
- Define tracepoints and write text queries
- Compile queries to advice
- Submit advice to PT agents

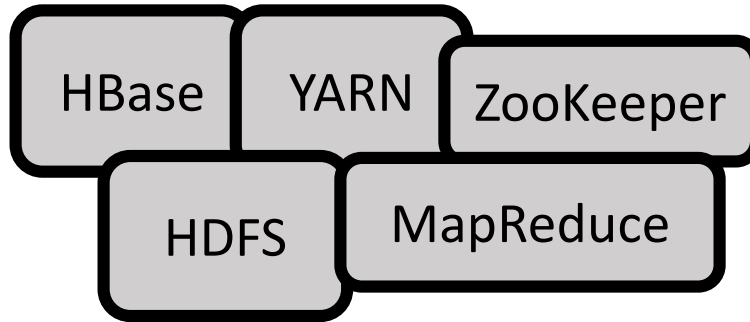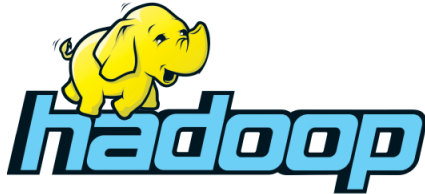Pivot Tracing Enabled    (+ 🧳 Baggage, + 🔧 PT Agent)

Pivot Tracing Enabled    (+ 🧳 Baggage, + 🔧 PT Agent)

HBase  YARN  ZooKeeper

HDFS  MapReduce

Adding Baggage: ~50-200 lines of code per system

Pivot Tracing Enabled (+ 🧳 Baggage, + 🔧 PT Agent)
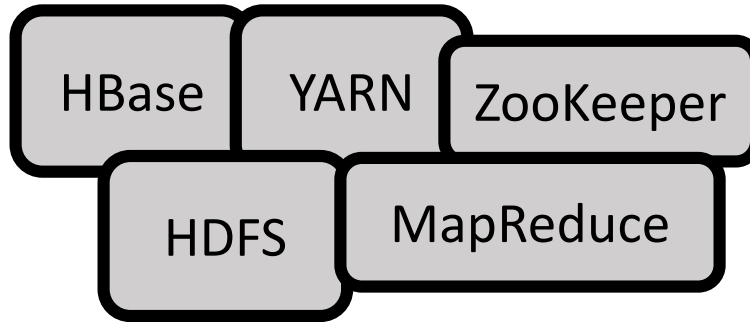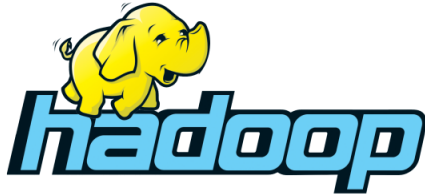
HBase  YARN  ZooKeeper

HDFS  MapReduce

Adding Baggage: ~50-200 lines of code per system

Primarily modifying execution boundaries:

Thread, Runnable, Callable, Queue

RPC invocations

Pivot Tracing Enabled    (+ 🧳 Baggage, + 🔧 PT Agent)



Adding Baggage: ~50-200 lines of code per system

Primarily modifying execution boundaries:
Thread, Runnable, Callable, Queue
RPC invocations

# Pivot Tracing Overheads

- **Pivot Tracing Enabled**  (+ 🧳 Baggage, + 🔑 PT Agent)

    Application level benchmarks: **baseline** 0.3% overhead

# Pivot Tracing Overheads

- **Pivot Tracing Enabled** (+ 🧳 Baggage, + 🔨 PT Agent)

   Application level benchmarks: **baseline** 0.3% overhead

- No overhead for queries / tracepoints until installed

# Pivot Tracing Overheads

- **Pivot Tracing Enabled** (+ 🧳 Baggage, + 🔑 PT Agent)

  Application level benchmarks: **baseline** 0.3% overhead

- No overhead for queries / tracepoints until installed

- With queries from paper installed

  Application level benchmarks: max 14.3% overhead

  (CPU-only lookups)

# Pivot Tracing Overheads

- ## Pivot Tracing Enabled (+ 🧳 Baggage, + 🔧 PT Agent)
  Application level benchmarks: **baseline** 0.3% overhead

- ## No overhead for queries / tracepoints until installed

- ## With queries from paper installed
  Application level benchmarks: max 14.3% overhead

  (CPU-only lookups)
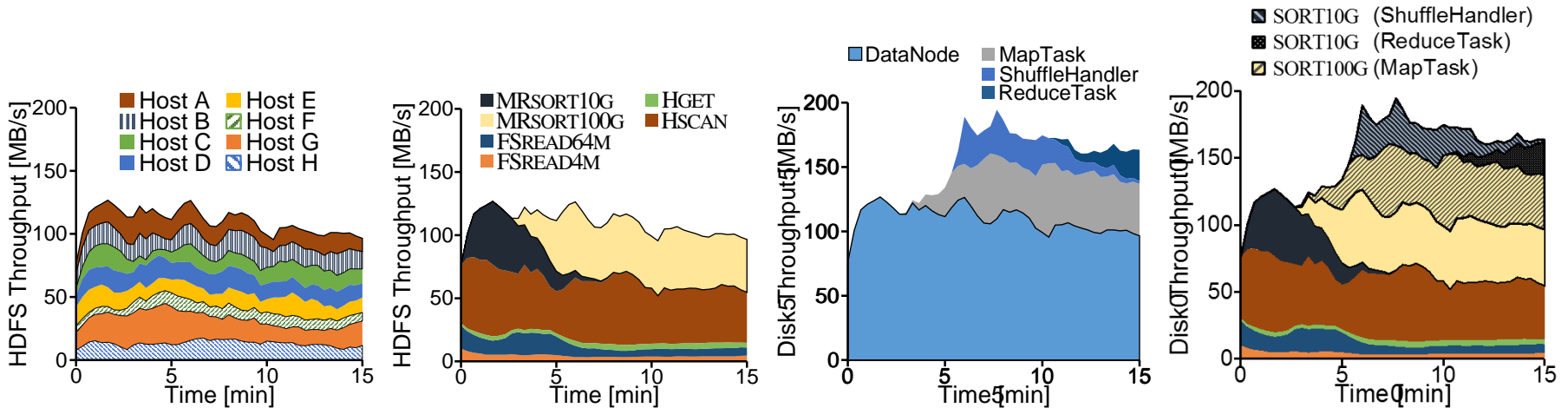
  Largest baggage size: ~137 bytes

# Experiments

34

# Experiments

## 1. Monitoring queries with various groupings

# Experiments

1. Monitoring queries with various groupings

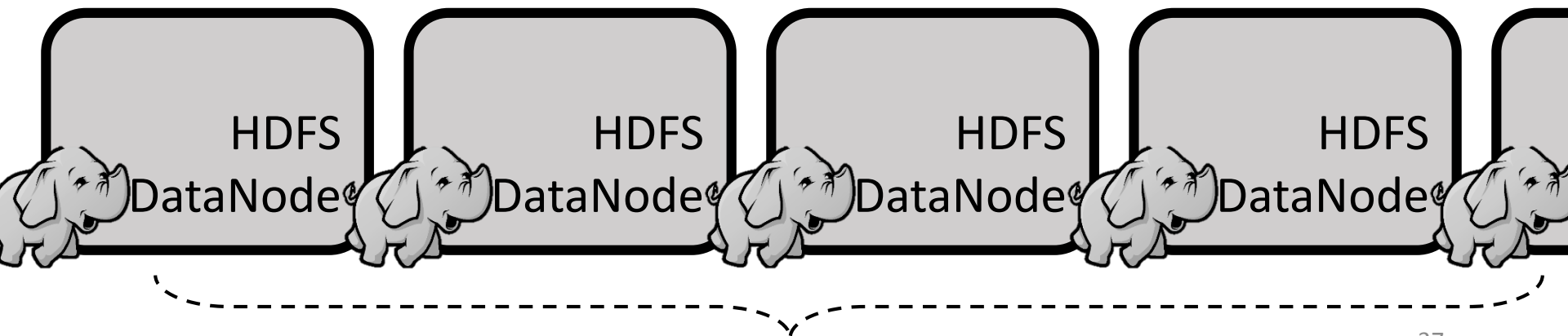2. Decomposing request latencies

# Experiments

1. Monitoring queries with various groupings

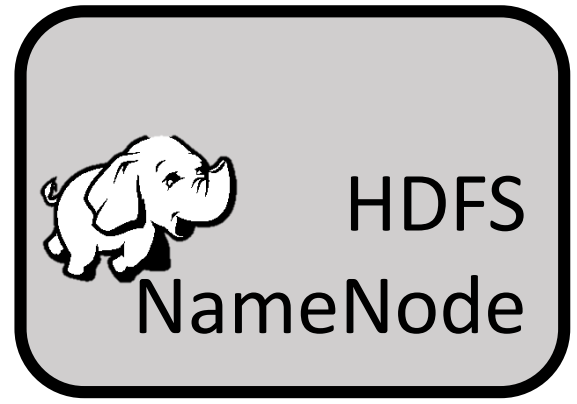2. Decomposing request latencies

3. Debugging recurring problems

# Experiments

1. Monitoring queries with various groupings

2. Decomposing request latencies
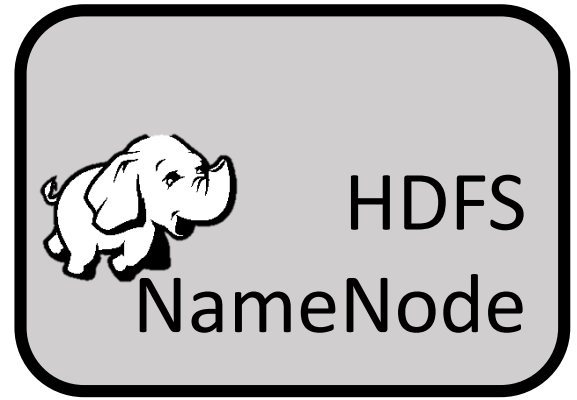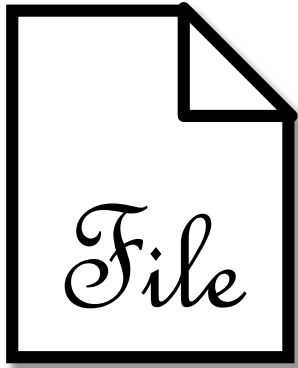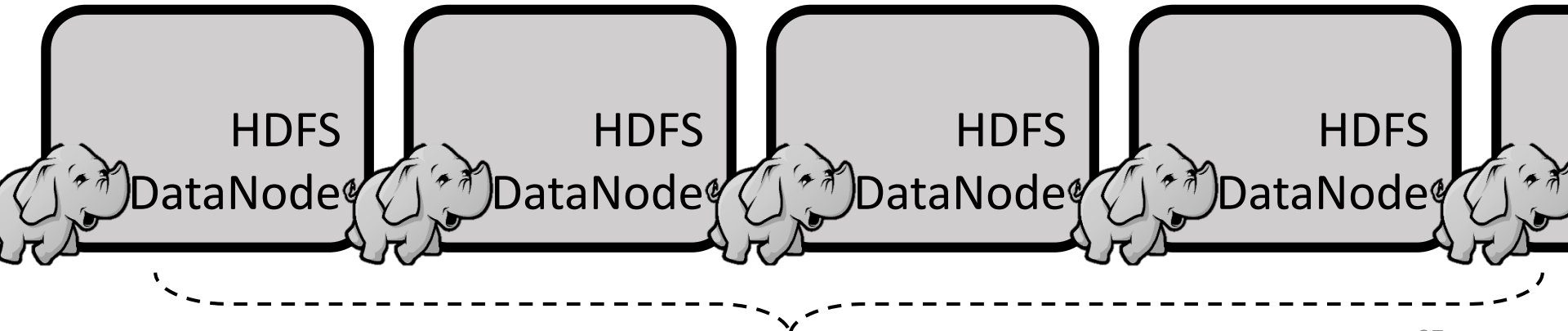
3. Debugging recurring problems

Filesystem Metadata

HDFS
NameNode

Filesystem Metadata

HDFS NameNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

Replicated block storage

File

HDFS NameNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

Replicated block storage

37

File

Filesystem Metadata

HDFS
NameNode

HDFS
DataNode

File

HDFS
DataNode

File

HDFS
DataNode

HDFS
DataNode

Replicated block storage

37

File

File 2 3 5

HDFS NameNode

File

File

HDFS DataNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

Replicated block storage

37

File 2 3 5

HDFS NameNode

File

File

HDFS DataNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

Replicated block storage

Client

File  2  3  5

Filesystem Metadata

HDFS NameNode

HDFS DataNode

File

HDFS DataNode

File

HDFS DataNode

HDFS DataNode

Replicated block storage

Client

File

Filesystem Metadata

File 2 3 5

GetBlockLocations

HDFS NameNode

File

File

HDFS DataNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

Replicated block storage

37

Client

File  2 3 5

GetBlockLocations

HDFS NameNode

Filesystem Metadata

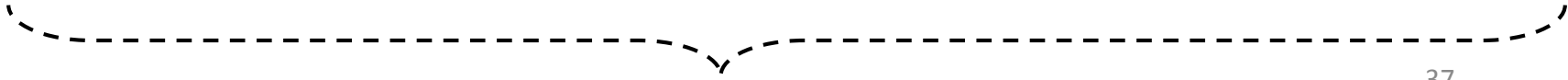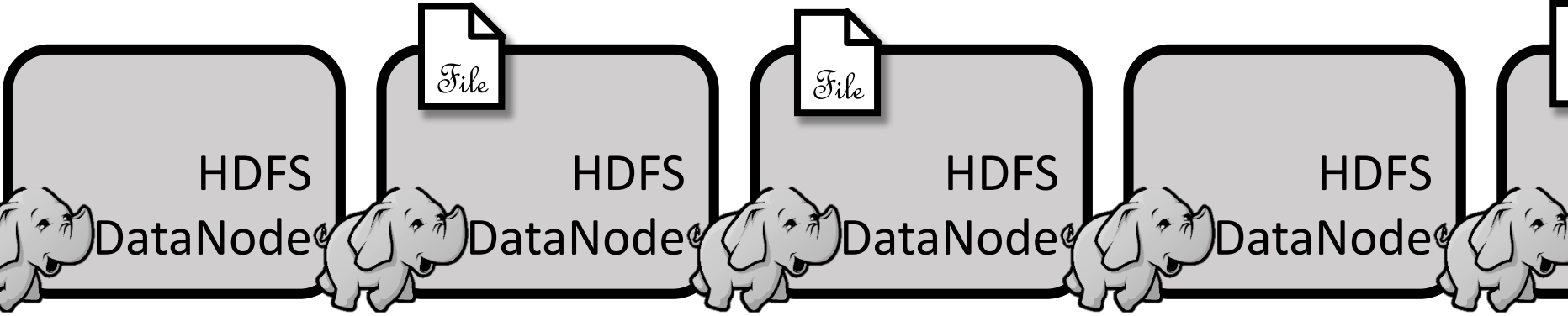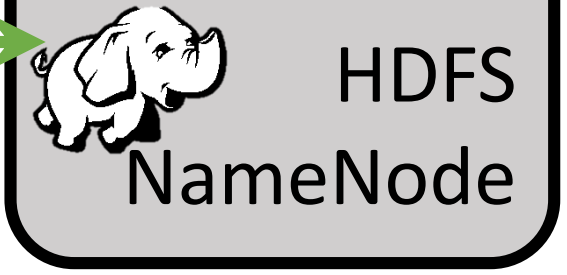File  2 3 5

File

File

HDFS DataNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

Replicated block storage

37

Filesystem Metadata

Client

File 2 3 5
GetBlockLocations

HDFS NameNode

File 2 3 5

**HDFS Replica Selection Policy**

Same machine?
-> read local

Same rack?
-> read rack-local

Otherwise
-> select randomly

File

File

HDFS DataNode

HDFS DataNode

HDFS DataNode

HDFS DataNode

Replicated block storage

37

Filesystem Metadata

File **2** **3** **5**
GetBlockLocations

HDFS NameNode

Client

HDFS Replica Selection Policy

Same machine?
-> read local

Same rack?
-> read rack-local

Otherwise
-> select randomly

File

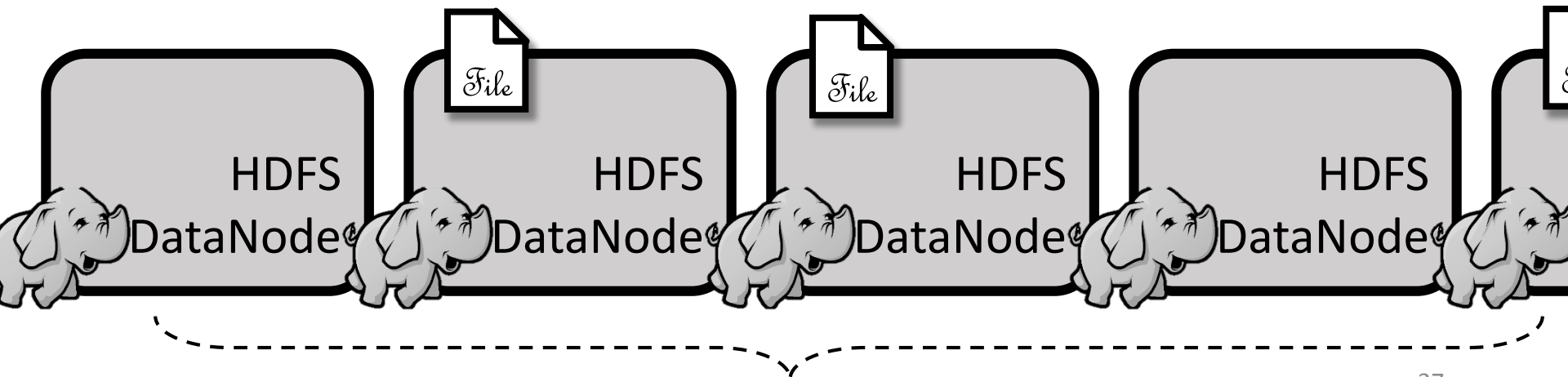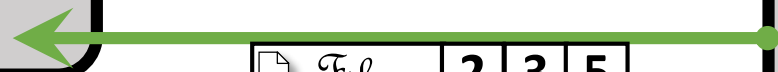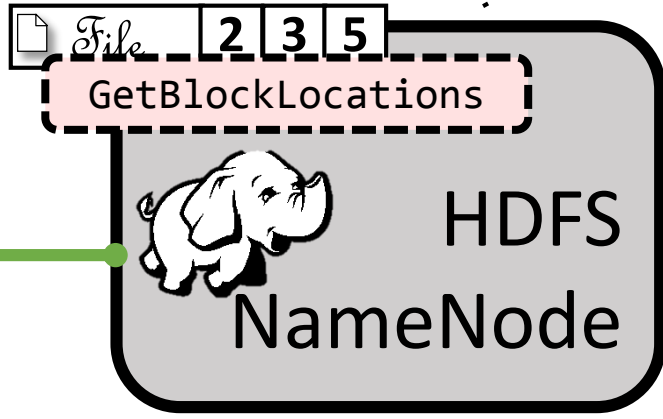HDFS DataNode

File

HDFS DataNode

File

DataTransferProtocol

HDFS DataNode
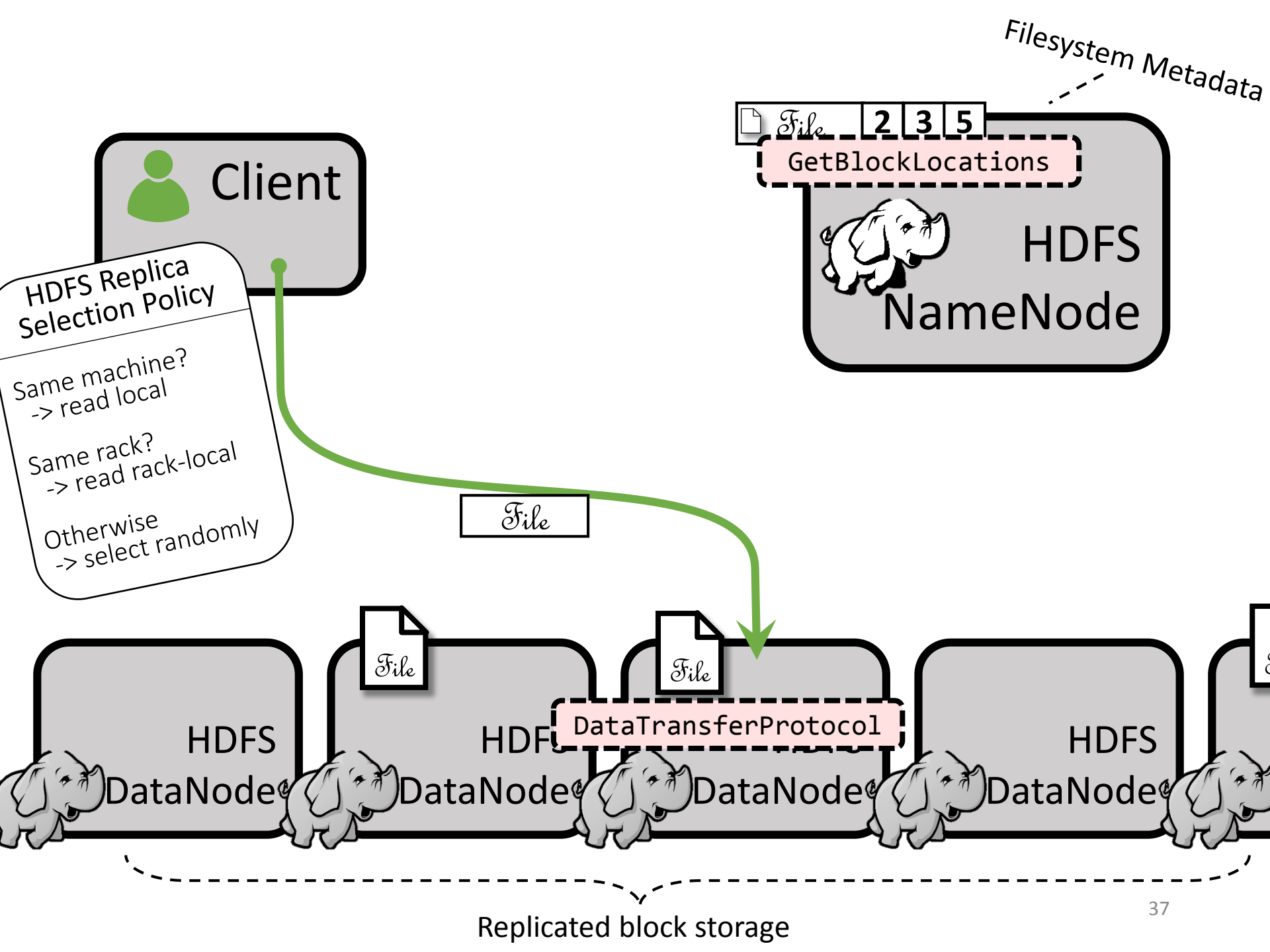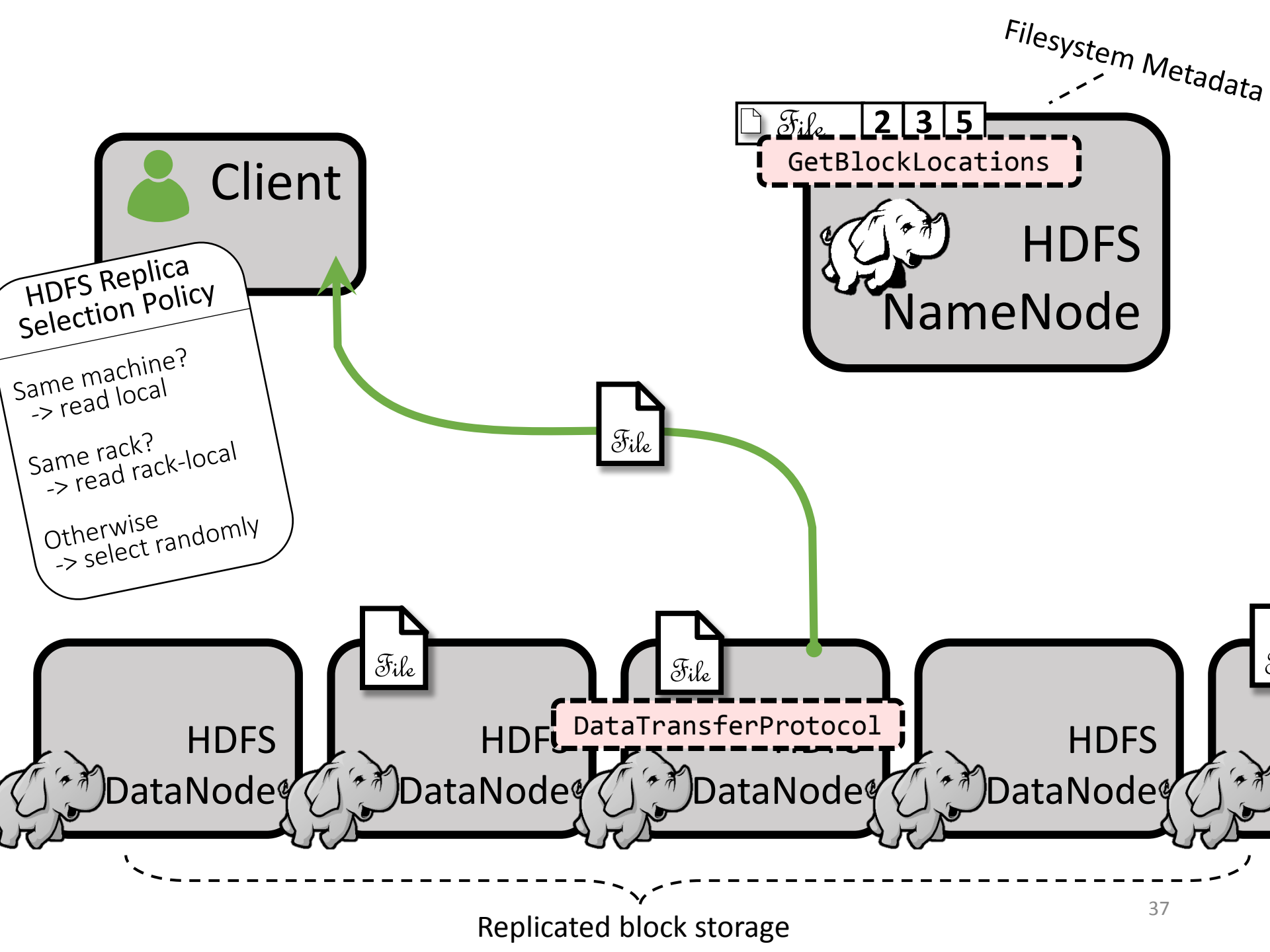
HDFS DataNode

Replicated block storage

Filesystem Metadata

Client

📄 𝓕𝓲𝓵𝓮  **2** **3** **5**

GetBlockLocations

HDFS NameNode

HDFS Replica Selection Policy

Same machine?
-> read local

Same rack?
-> read rack-local

Otherwise
-> select randomly

𝓕𝓲𝓵𝓮

𝓕𝓲𝓵𝓮

𝓕𝓲𝓵𝓮

HDFS DataNode

HDFS DataNode

HDFS DataNode

DataTransferProtocol

HDFS DataNode

Replicated block storage

37

**Host A** **Host B** **Host C** **Host D**

**Host E** **Host F** **Host G** **Host H**

**8 Worker Hosts**

Client Workload Generator
- Randomly read from large dataset

HDFS DataNode

**Host A**   **Host B**   **Host C**   **Host D**

**Host E**   **Host F**   **Host G**   **Host H**

+ HDFS NameNode

**8 Worker Hosts**

Client Workload Generator
- Randomly read from large dataset

HDFS DataNode

# Same machines, same processes, same workloads

# Same machines, same processes, same workloads

I expected uniform throughput from workload generators

# Same machines, same processes, same workloads

## I expected uniform throughput from workload generators



Consistently lower throughput

# Same machines, same processes, same workloads

I expected uniform throughput from workload generators
I expected uniform throughput on DataNodes



Consistently lower throughput

# Same machines, same processes, same workloads

I expected uniform throughput from workload generators
I expected uniform throughput on DataNodes



Large variations in DataNode throughput

It's probably a bug in the workload generator I wrote

My hypothesis:

Workload generator is not randomly looking up files

It's probably a bug in the workload generator I wrote

My hypothesis:

Workload generator is not randomly looking up files

```
From blockLocations In NameNode.GetBlockLocations
  GroupBy blockLocations.fileName
  Select blockLocations.fileName, COUNT
```

GetBlockLocations

HDFS
NameNode

# It's probably a bug in the workload generator I wrote

## My hypothesis:

Workload generator is not randomly looking up files

```
From blockLocations In NameNode.GetBlockLocations
  GroupBy blockLocations.fileName
  Select blockLocations.fileName, COUNT
```



GetBlockLocations

HDFS
NameNode

# It's probably a bug in the workload generator I wrote

## My hypothesis:

Workload generator is not randomly looking up files

```
From blockLocations In NameNode.GetBlockLocations
  Join cl In Client.DoRandomRead On cl -> blockLocations
  GroupBy cl.host, blockLocations.fileName
  Select cl.host, blockLocations.fileName, COUNT
```



Client

DoRandomRead

GetBlockLocations

HDFS
NameNode

# It's probably a bug in the workload generator I wrote

## My hypothesis:

Workload generator is not randomly looking up files

```
From blockLocations In NameNode.GetBlockLocations
    Join cl In Client.DoRandomRead On cl -> blockLocations
GroupBy cl.host, blockLocations.fileName
Select cl.host, blockLocations.fileName, COUNT
```



Client

DoRandomRead

GetBlockLocations

HDFS
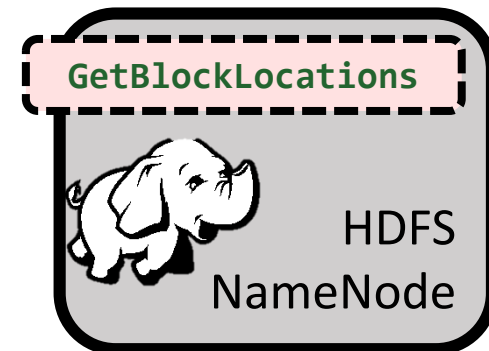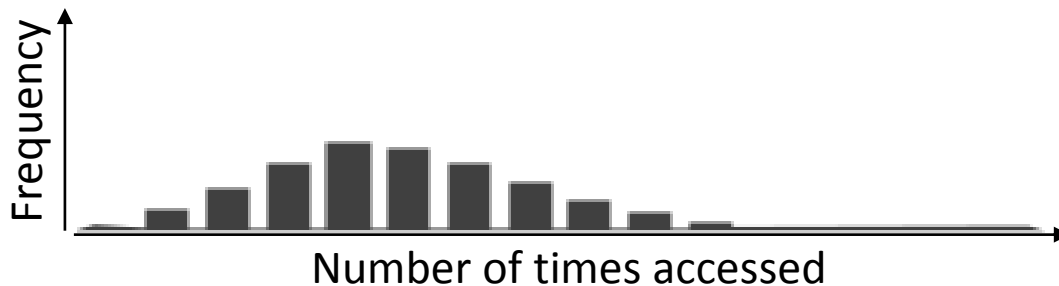NameNode
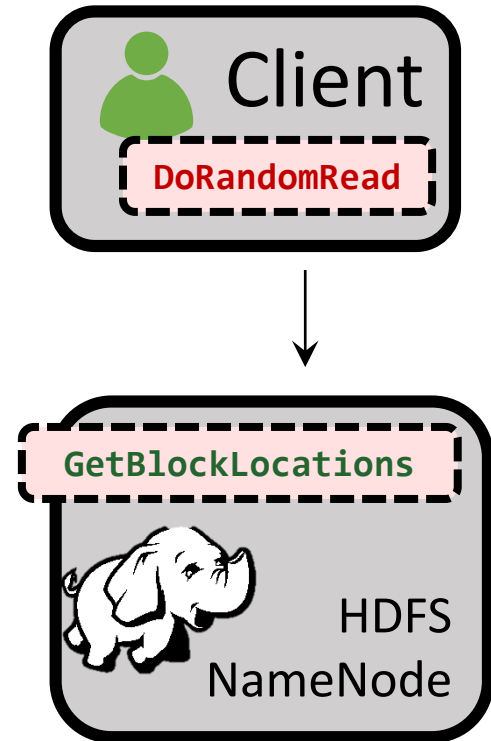
Frequency

Number of times accessed

# It's probably a bug in the workload generator I wrote

## My hypothesis:

Workload generator is not randomly looking up files

```
From blockLocations In NameNode.GetBlockLocations
  Join cl In Client.DoRandomRead On cl -> blockLocations
GroupBy cl.host, blockLocations.fileName
Select cl.host, blockLocations.fileName, COUNT
```



Client

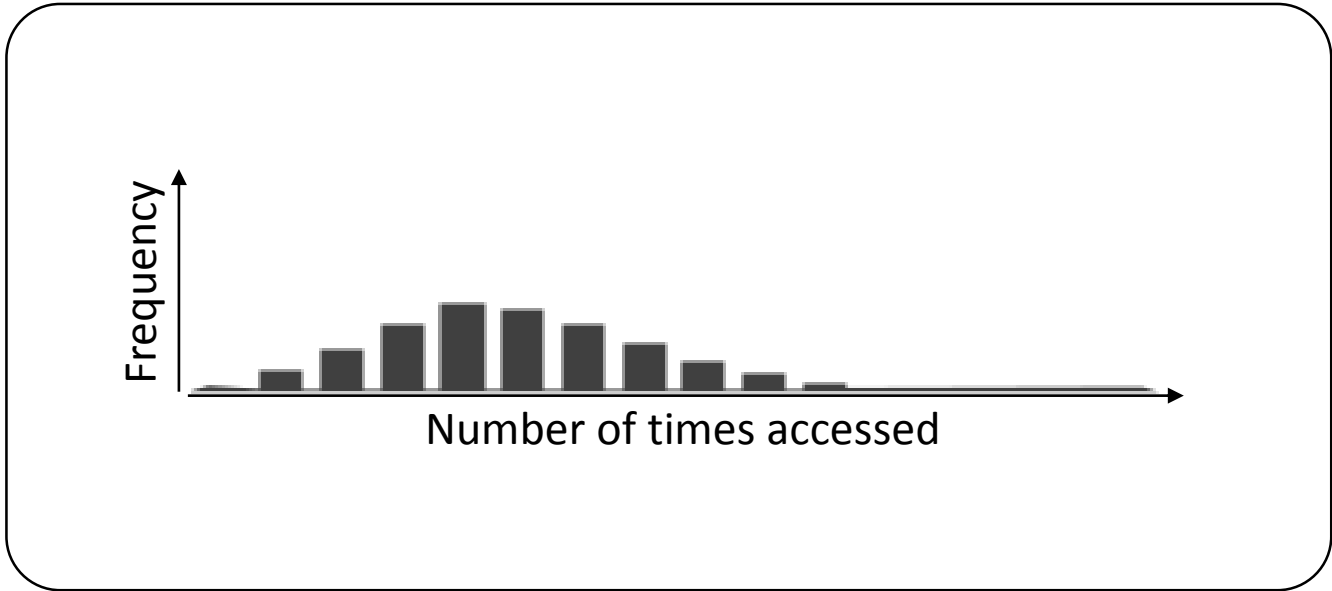**DoRandomRead**

**GetBlockLocations**

HDFS
NameNode

# It's probably a bug in the workload generator I wrote

## My hypothesis:

Workload generator is not randomly looking up files

```
From blockLocations In NameNode.GetBlockLocations
  Join cl In Client.DoRandomRead On cl -> blockLocations
  GroupBy cl.host, blockLocations.fileName
Select cl.host, blockLocations.fileName, COUNT
```
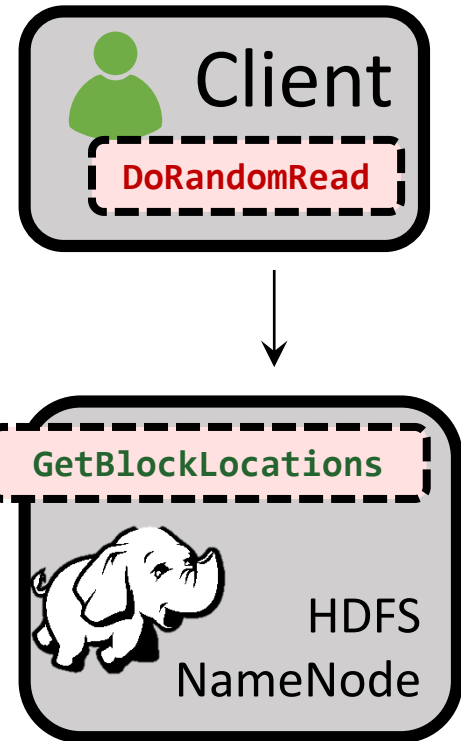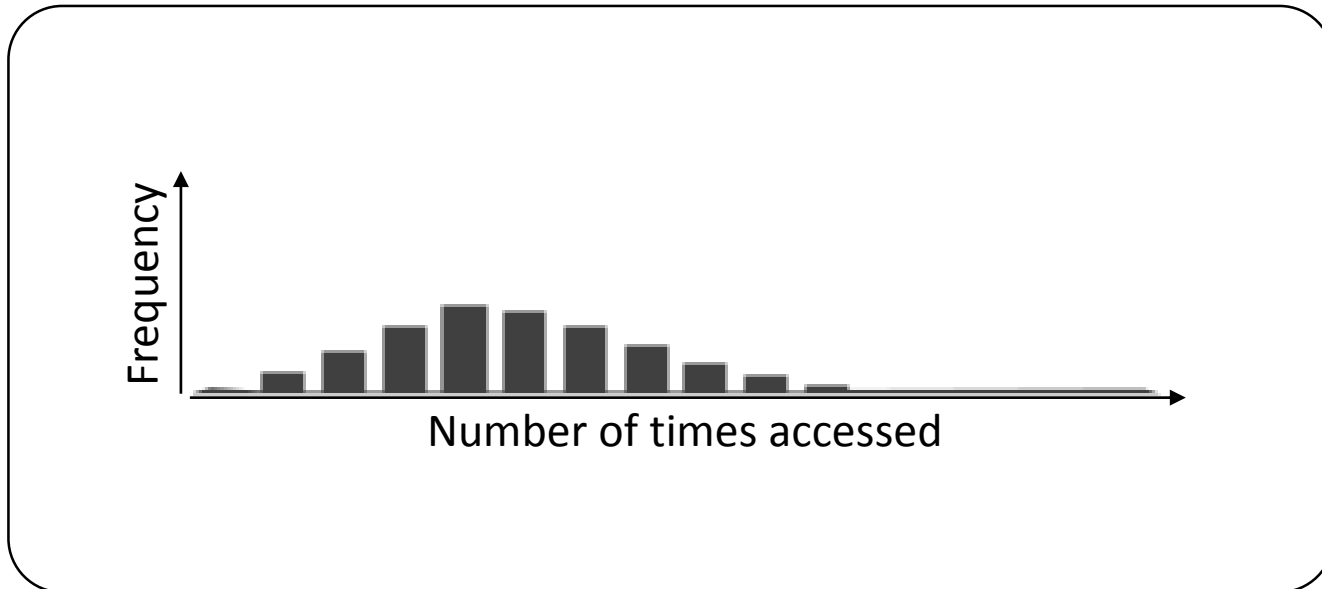


Client

DoRandomRead

GetBlockLocations

HDFS
NameNode

Frequency

Number of times accessed

# It's probably a bug in the workload generator I wrote

## My hypothesis:

Workload generator is not randomly looking up files

```
From blockLocations In NameNode.GetBlockLocations
  Join cl In Client.DoRandomRead On cl -> blockLocations
  GroupBy cl.host, blockLocations.fileName
  Select cl.host, blockLocations.fileName, COUNT
```
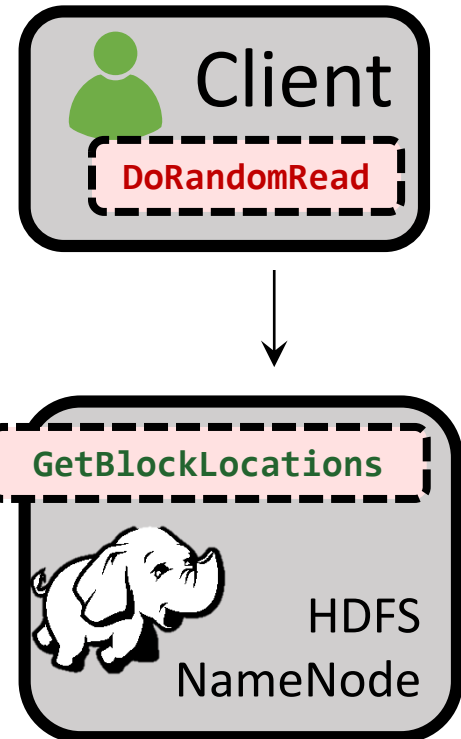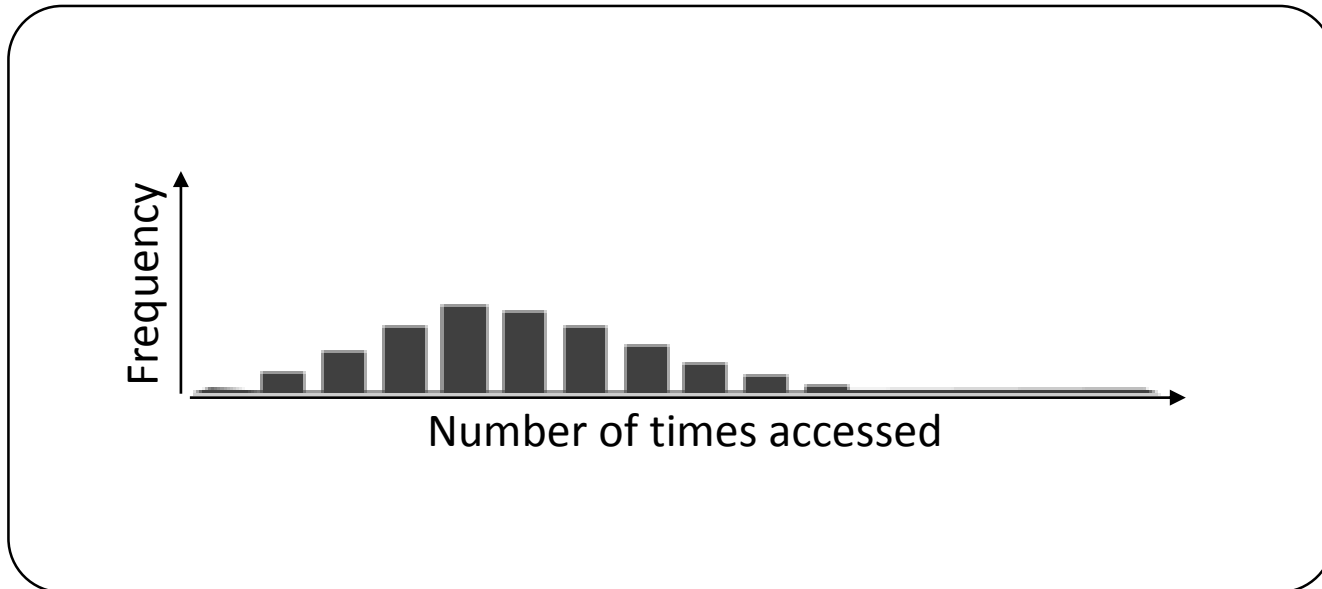
# It's probably a bug in the workload generator I wrote

## My hypothesis:

Workload generator is not randomly looking up files

```
From blockLocations In NameNode.GetBlockLocations
   Join cl In Client.DoRandomRead On cl -> blockLocations
   GroupBy cl.host, blockLocations.fileName
   Select cl.host, blockLocations.fileName, COUNT
```
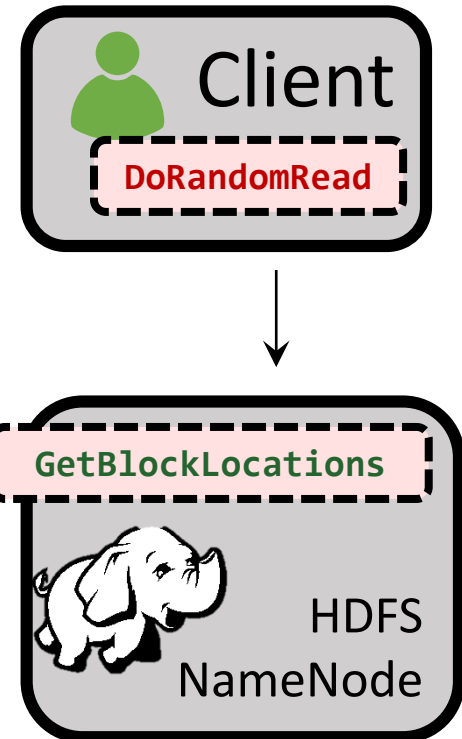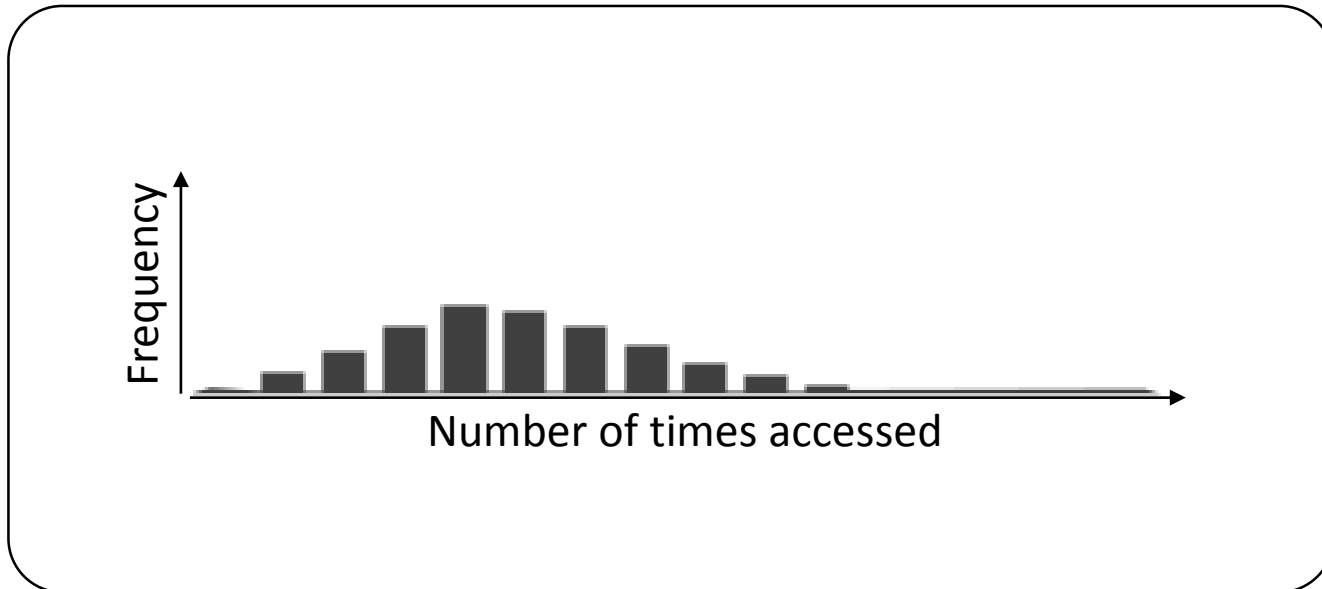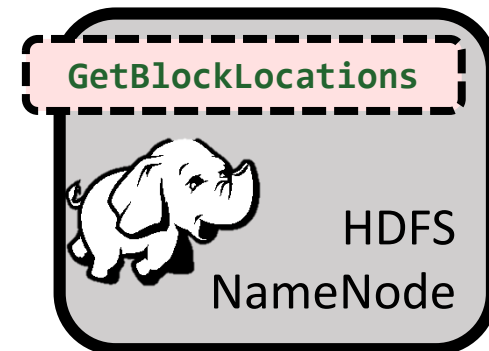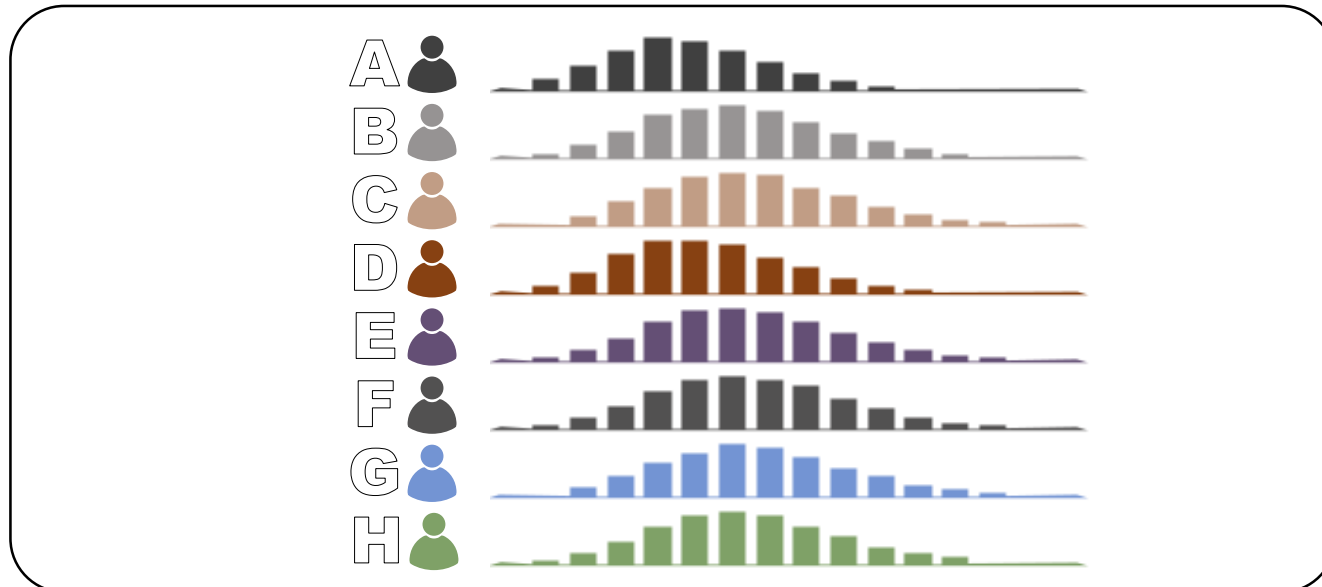
Maybe skewed DataNode throughput is because some DataNodes store more files than others
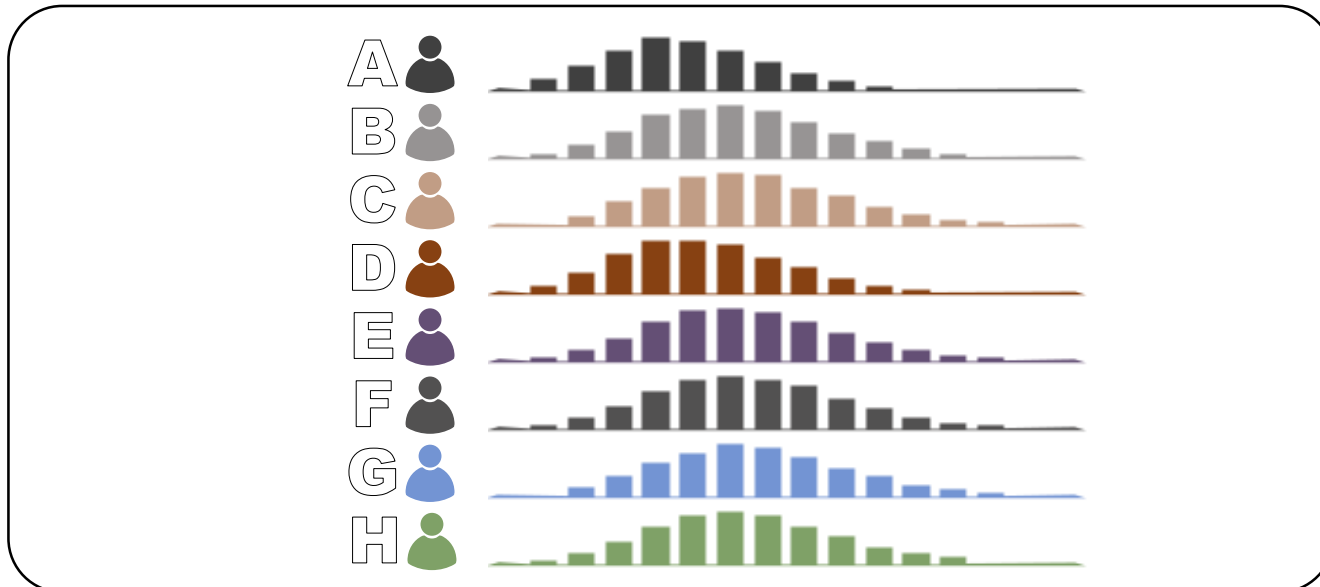
# Maybe skewed DataNode throughput is because some DataNodes store more files than others



How often was each DataNode a replica host?

# Maybe skewed DataNode throughput is because some DataNodes store more files than others

How often was each DataNode a replica host?

```
From blockLocations In NameNode.GetBlockLocations
  GroupBy blockLocations.replicas
  Select blockLocations.replicas, COUNT
```

GetBlockLocations

HDFS
NameNode

# Maybe skewed DataNode throughput is because some DataNodes store more files than others

How often was each DataNode a replica host?

```
From blockLocations In NameNode.GetBlockLocations
  GroupBy blockLocations.replicas
  Select blockLocations.replicas, COUNT
```



Replica Location



GetBlockLocations

HDFS NameNode

# Maybe skewed DataNode throughput is because some DataNodes store more files than others



How often was each DataNode a replica host?

```
From blockLocations In NameNode.GetBlockLocations
  Join cl In Client.DoRandomRead On cl -> blockLocations
  GroupBy cl.host, blockLocations.replicas
  Select cl.host, blockLocations.replicas, COUNT
```
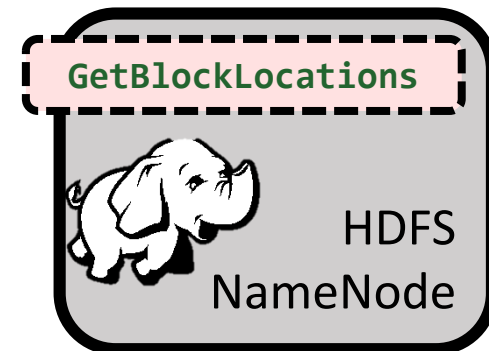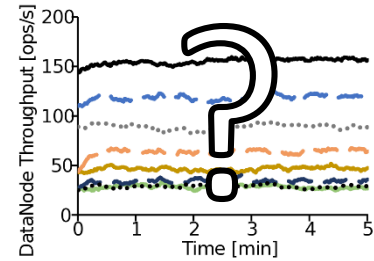


Count

A B C D E F G H

Replica Location
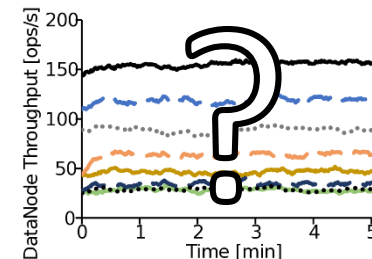
Client

**DoRandomRead**

**GetBlockLocations**

HDFS
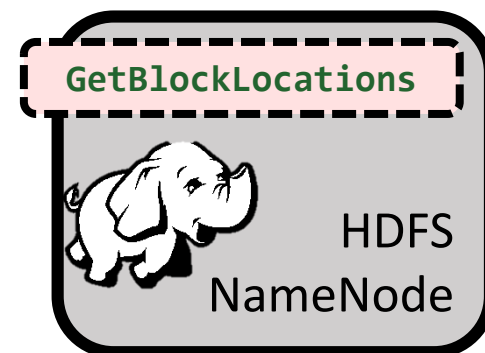NameNode

# Maybe skewed DataNode throughput is because some DataNodes store more files than others
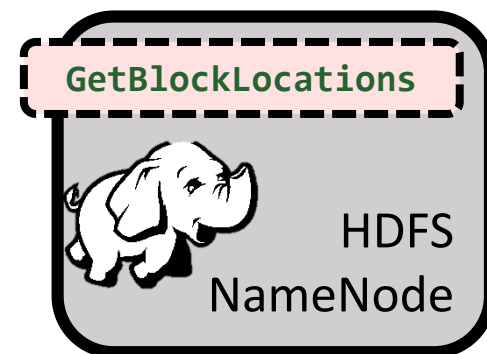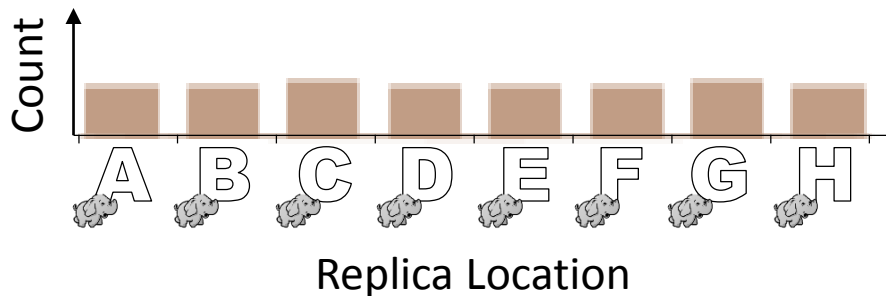


How often was each DataNode a replica host?

```
From blockLocations In NameNode.GetBlockLocations
   Join cl In Client.DoRandomRead On cl -> blockLocations
   GroupBy cl.host, blockLocations.replicas
   Select cl.host, blockLocations.replicas, COUNT
```



Client

Replica Location



Client

**DoRandomRead**

**GetBlockLocations**

HDFS
NameNode

# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random

# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random

# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random



Hypothesis: choice of replica isn't random?

## Conclusions so far:

Clients are selecting files uniformly at random

Files are distributed across DNs uniformly at random

## Hypothesis: choice of replica isn't random?



HDFS Replica
Selection Policy

Same machine?
-> read local

Same rack?
-> read rack-local

Otherwise
-> select randomly

# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random



# Hypothesis: choice of replica isn't random?

When a file is read from a DataNode,
where else *could* it have been read from?

**HDFS Replica Selection Policy**

Same machine?
-> read local

Same rack?
-> read rack-local

Otherwise
-> select randomly

# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random



## Hypothesis: choice of replica isn't random?

When a file is read from a DataNode,
where else *could* it have been read from?

**HDFS Replica Selection Policy**

Same machine?
-> read local

Same rack?
-> read rack-local

Otherwise
-> select randomly

`GetBlockLocations`

HDFS
NameNode

`DataTransferProtocol`

HDFS
DataNode

42

# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random

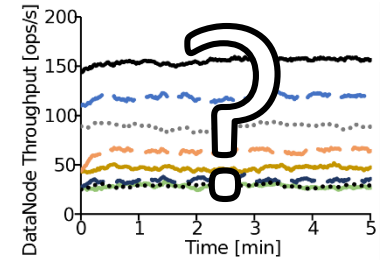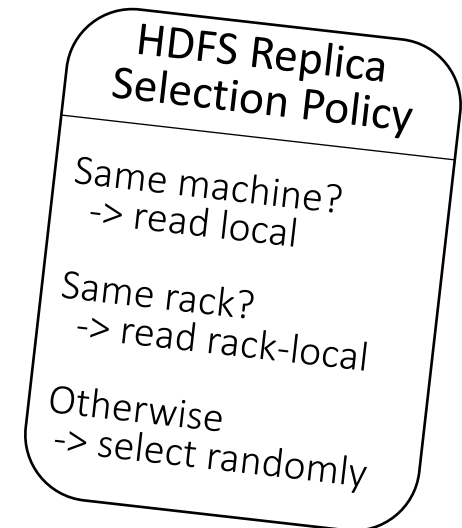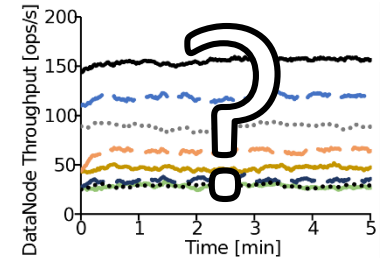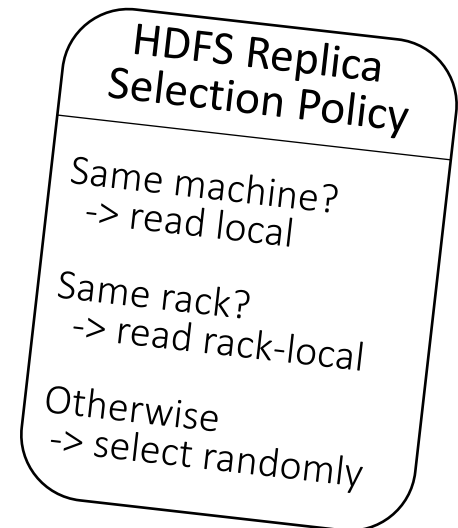# Hypothesis: choice of replica isn't random?

When a file is read from a DataNode,
where else *could* it have been read from?

**HDFS Replica Selection Policy**

Same machine?
-> read local

Same rack?
-> read rack-local

Otherwise
-> select randomly

**GetBlockLocations**

HDFS NameNode

**DataTransferProtocol**

HDFS DataNode

```
From readBlock In DataNode.DataTransferProtocol
    Join blockLocations In NameNode.GetBlockLocations
                    On blockLocations -> readBlock
    GroupBy blockLocations.replicas, readBlock.host
    Select blockLocations.replicas, readBlock.host, COUNT
```

42

# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random



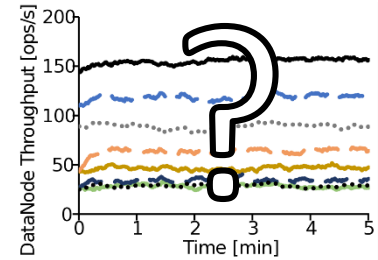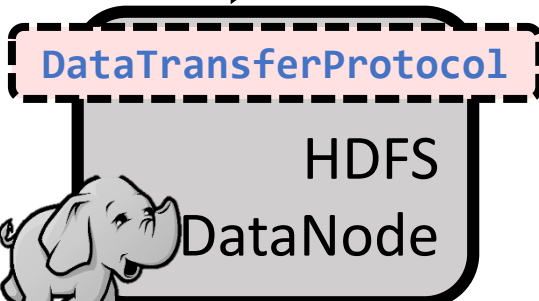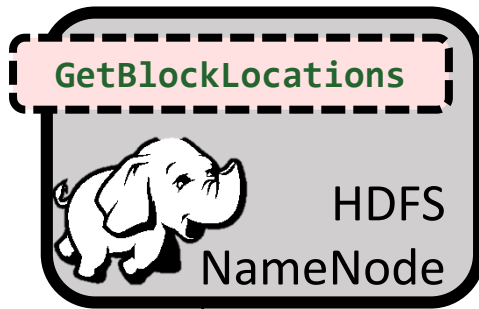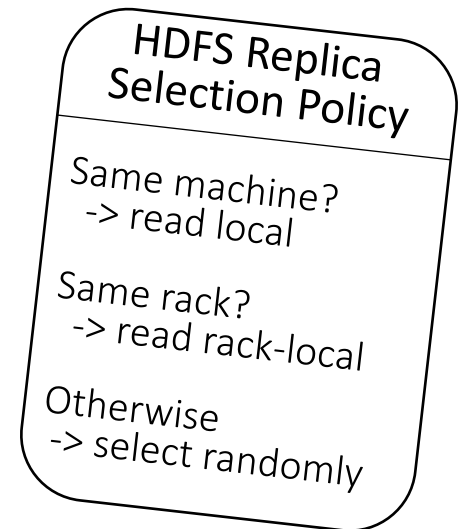# Hypothesis: choice of replica isn't random?
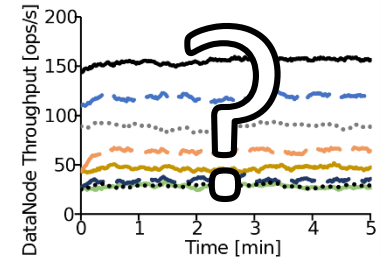
When a file is read from a DataNode,
where else *could* it have been read from?

**HDFS Replica Selection Policy**

~~Same machine? -> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

**GetBlockLocations**

HDFS NameNode

**DataTransferProtocol**

HDFS DataNode

```
From readBlock In DataNode.DataTransferProtocol
    Join blockLocations In NameNode.GetBlockLocations
                      On blockLocations -> readBlock
    GroupBy blockLocations.replicas, readBlock.host
    Select blockLocations.replicas, readBlock.host, COUNT
```
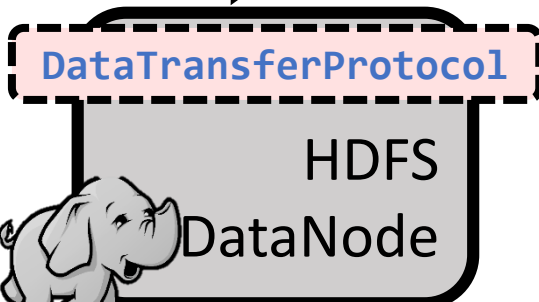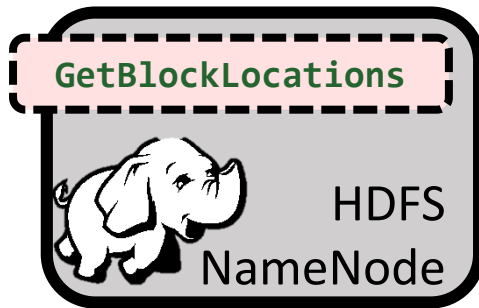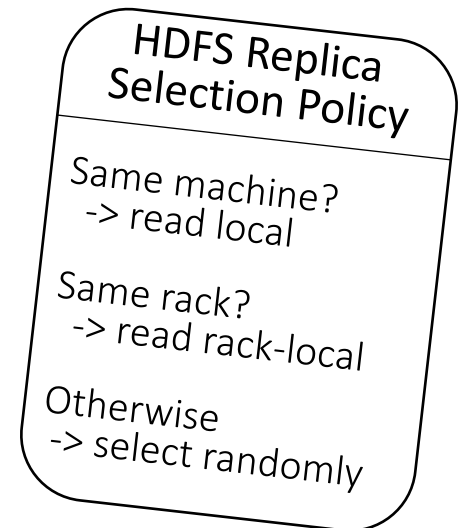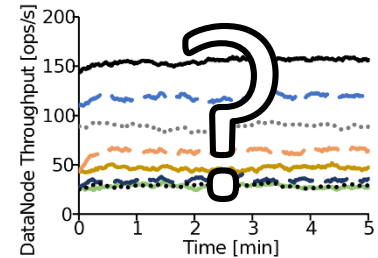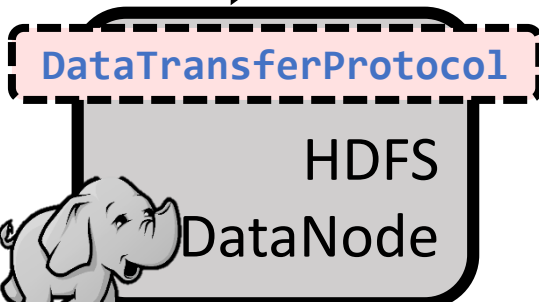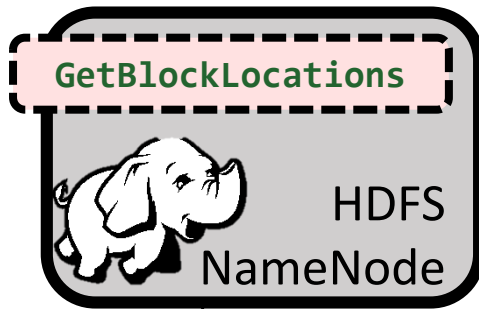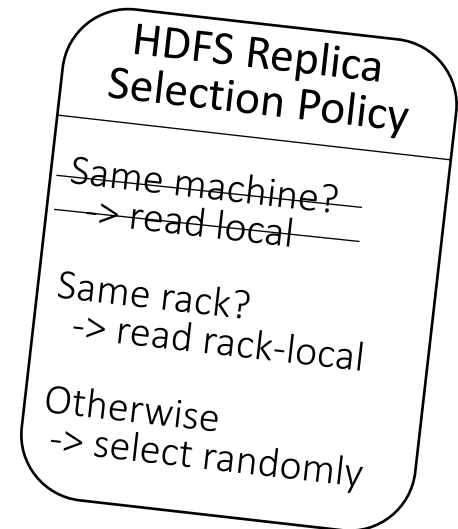
# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random



# Hypothesis: choice of replica isn't random?

**Client**

`DoRandomRead`

When a file is read from a DataNode,
where else *could* it have been read from?

**HDFS Replica Selection Policy**

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

`GetBlockLocations`

HDFS
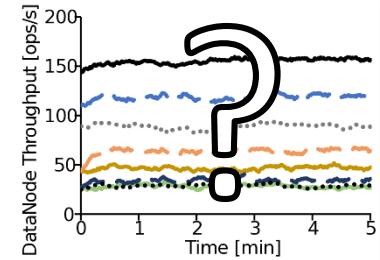NameNode

```
From readBlock In DataNode.DataTransferProtocol
    Join blockLocations In NameNode.GetBlockLocations
                On blockLocations -> readBlock
    GroupBy blockLocations.replicas, readBlock.host
    Select blockLocations.replicas, readBlock.host, COUNT
```
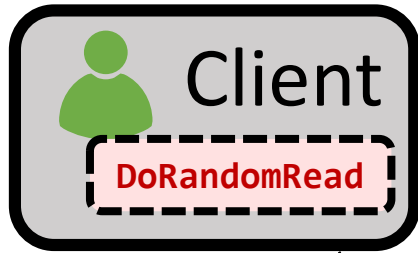
`DataTransferProtocol`

HDFS
DataNode

# Conclusions so far:

Clients are selecting files uniformly at random
Files are distributed across DNs uniformly at random



## Hypothesis: choice of replica isn't random?

**Client**

**DoRandomRead**

When a file is read from a DataNode,
where else *could* it have been read from?

**HDFS Replica Selection Policy**

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

**GetBlockLocations**

HDFS NameNode

**DataTransferProtocol**

HDFS DataNode

```
From readBlock In DataNode.DataTransferProtocol
    Join blockLocations In NameNode.GetBlockLocations
                    On blockLocations -> readBlock
    Join cl In Client.DoRandomRead
                    On cl -> blockLocations
    Where cl.host != readBlock.host
    GroupBy blockLocations.replicas, readBlock.host
    Select blockLocations.replicas, readBlock.host, COUNT
```

# Conclusions so far:

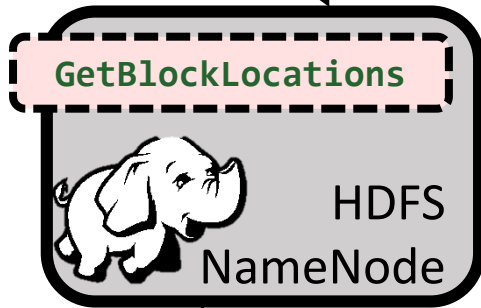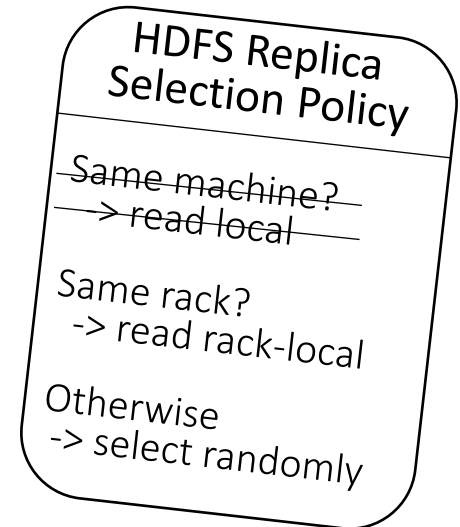Clients are selecting files uniformly at random
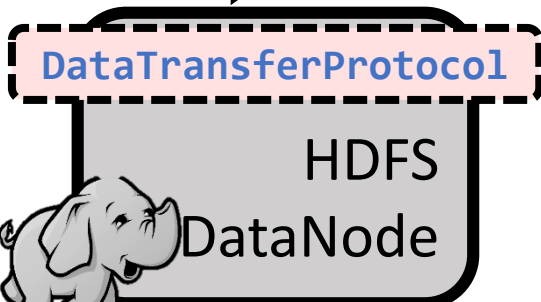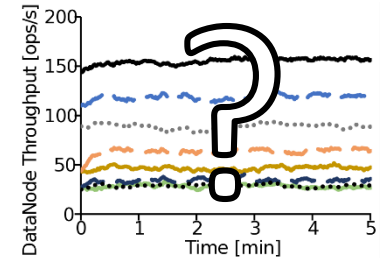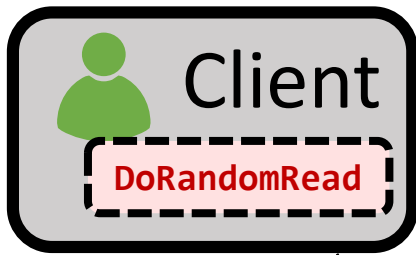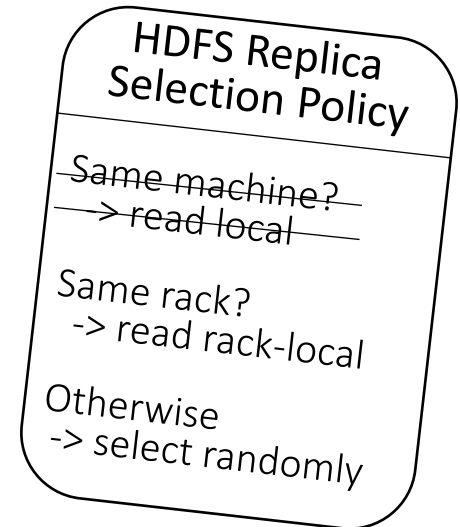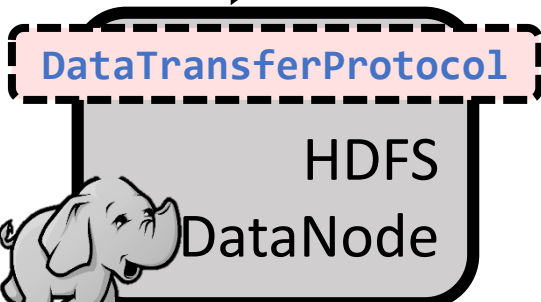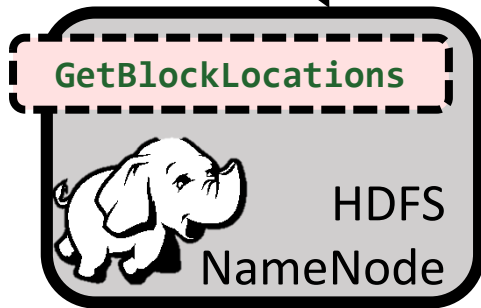Files are distributed across DNs uniformly at random



## Hypothesis: choice of replica isn't random?



When a file is read from a DataNode,
where else *could* it have been read from?

**HDFS Replica Selection Policy**

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly
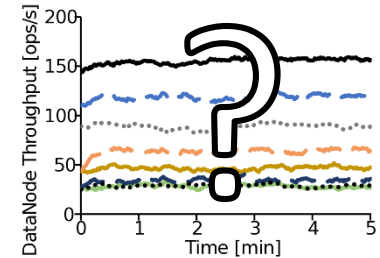
```
From readBlock In DataNode.DataTransferProtocol
    Join blockLocations In NameNode.GetBlockLocations
                On blockLocations -> readBlock
    Join cl In Client.DoRandomRead
                On cl -> blockLocations
    Where cl.host != readBlock.host
    GroupBy blockLocations.replicas, readBlock.host
    Select blockLocations.replicas, readBlock.host, COUNT
```
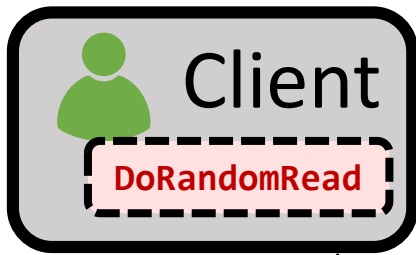
42

DataNode Throughput [ops/s]

200
150
100
50
0

Time [min]
0 1 2 3 4 5

**?**

Client

**DoRandomRead**

**GetBlockLocations**

HDFS
NameNode

**DataTransferProtocol**

HDFS
DataNode

HDFS Replica
Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

```
From readBlock In DataNode.DataTransferProtocol
   Join blockLocations In NameNode.GetBlockLocations
                 On blockLocations -> readBlock
   Join cl In Client.DoRandomRead
                 On cl -> blockLocations
   Where cl.host != readBlock.host
   GroupBy blockLocations.replicas, readBlock.host
   Select blockLocations.replicas, readBlock.host, COUNT
```

42

Client

**DoRandomRead**

**GetBlockLocations**

HDFS NameNode

**DataTransferProtocol**

HDFS DataNode

When both 𝐵 and 𝐷 host replicas,
Clients choose 𝐵 this often: ▬ (~50%)
Clients choose 𝐷 this often: ▬ (~50%)

HDFS Replica Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
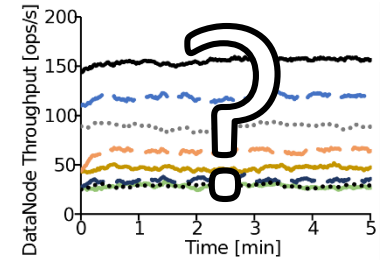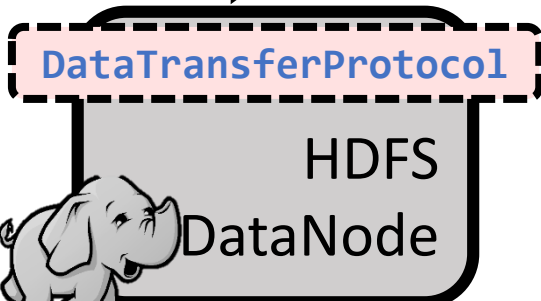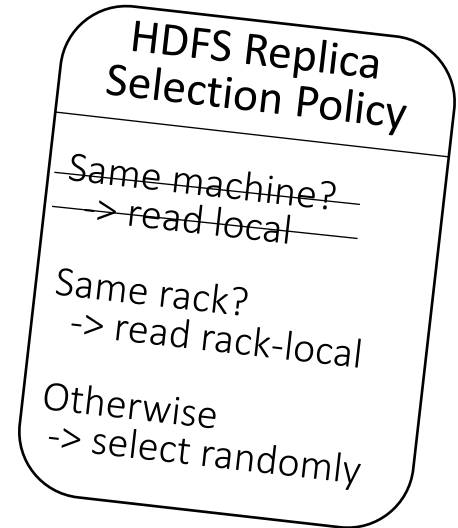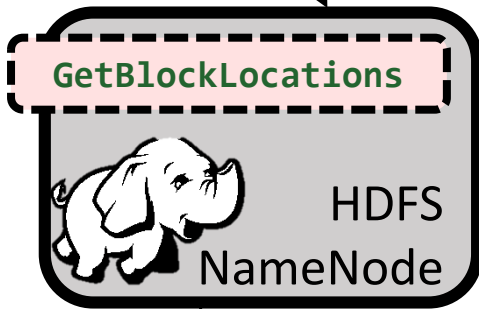-> select randomly

```
From readBlock In DataNode.DataTransferProtocol
   Join blockLocations In NameNode.GetBlockLocations
                  On blockLocations -> readBlock
   Join cl In Client.DoRandomRead
                  On cl -> blockLocations
   Where cl.host != readBlock.host
   GroupBy blockLocations.replicas, readBlock.host
   Select blockLocations.replicas, readBlock.host, COUNT
```

43

**Client**

DoRandomRead

GetBlockLocations

HDFS NameNode

DataTransferProtocol

HDFS DataNode

When both 🐘A and 🐘C host replicas, Clients choose 🐘A this often: ■ (100%)

HDFS Replica Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

```
From readBlock In DataNode.DataTransferProtocol
   Join blockLocations In NameNode.GetBlockLocations
                On blockLocations -> readBlock
   Join cl In Client.DoRandomRead
                On cl -> blockLocations
   Where cl.host != readBlock.host
   GroupBy blockLocations.replicas, readBlock.host
   Select blockLocations.replicas, readBlock.host, COUNT
```
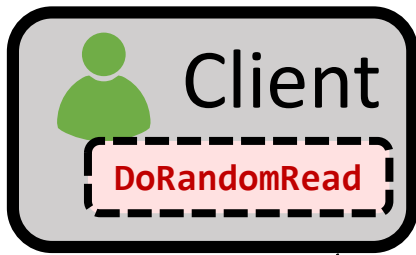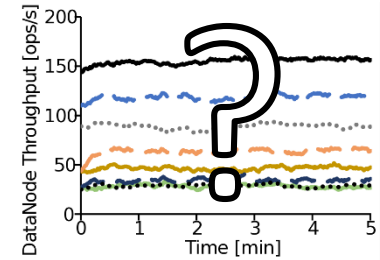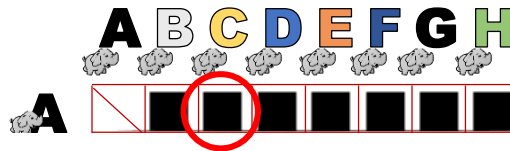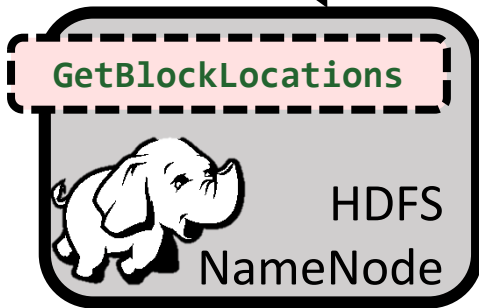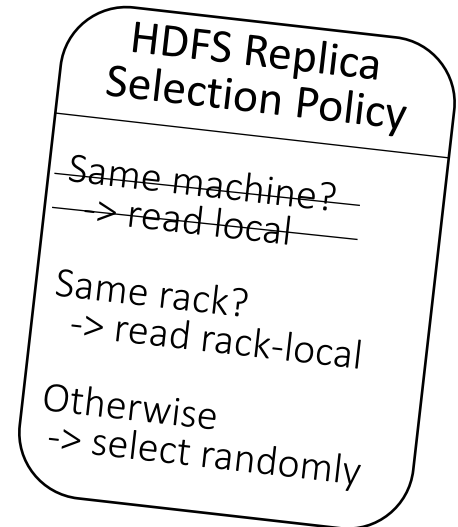
43

A B C D E F G H

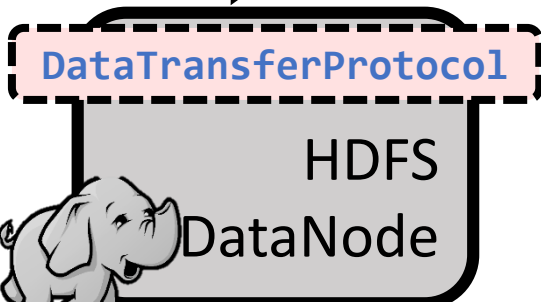DataNode Throughput [ops/s]

200
150
100
50
0

0 1 2 3 4 5
Time [min]

When both **A** and **C** host replicas,
Clients choose **A** this often: ■ (100%)
Clients choose **C** this often: __ (0%)

HDFS Replica Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

## Client

**DoRandomRead**

## HDFS NameNode

**GetBlockLocations**

## HDFS DataNode

**DataTransferProtocol**

```
From readBlock In DataNode.DataTransferProtocol
   Join blockLocations In NameNode.GetBlockLocations
              On blockLocations -> readBlock
   Join cl In Client.DoRandomRead
              On cl -> blockLocations
   Where cl.host != readBlock.host
   GroupBy blockLocations.replicas, readBlock.host
   Select blockLocations.replicas, readBlock.host, COUNT
```
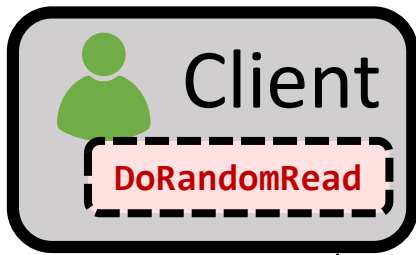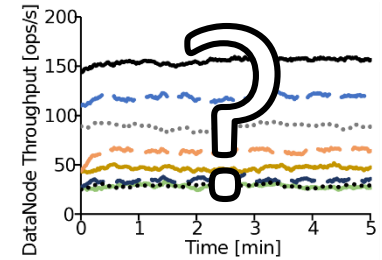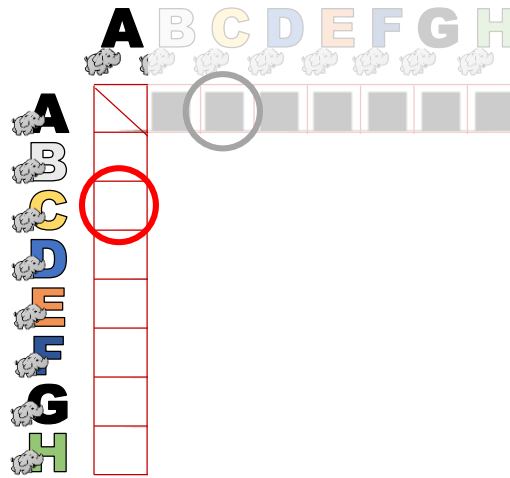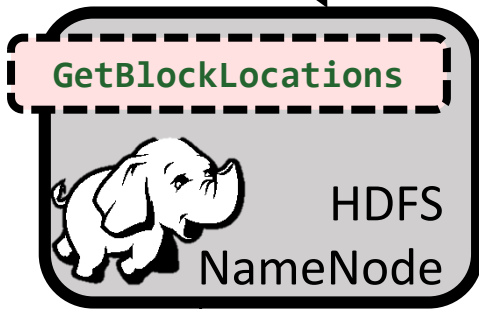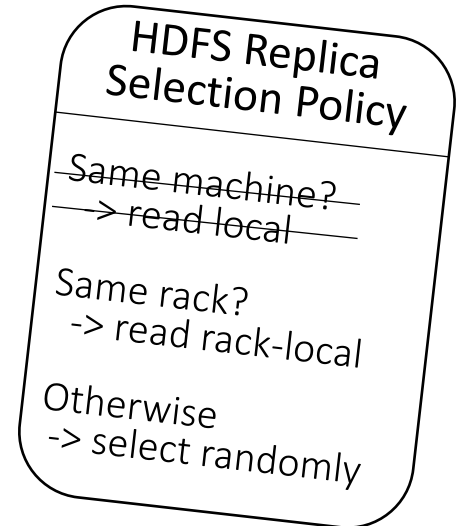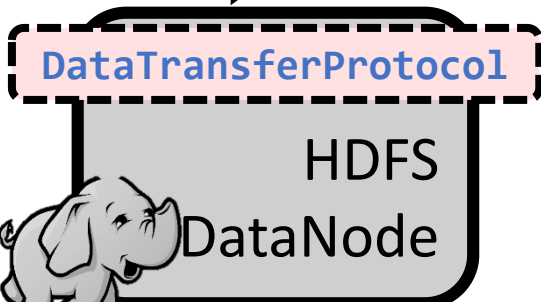
 A B C D E F G H

DataNode Throughput [ops/s]

When both 🐘A and 🐘C host replicas,
Clients choose 🐘A this often: ■ (100%)
Clients choose 🐘C this often: __ (0%)

HDFS Replica
Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

Client

**DoRandomRead**

**GetBlockLocations**

HDFS
NameNode
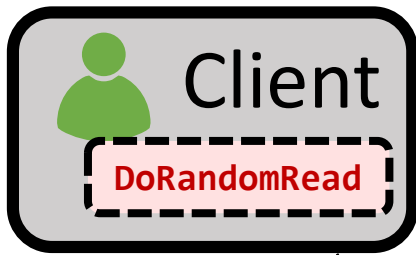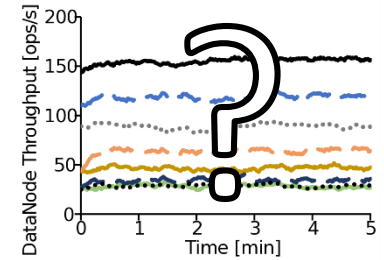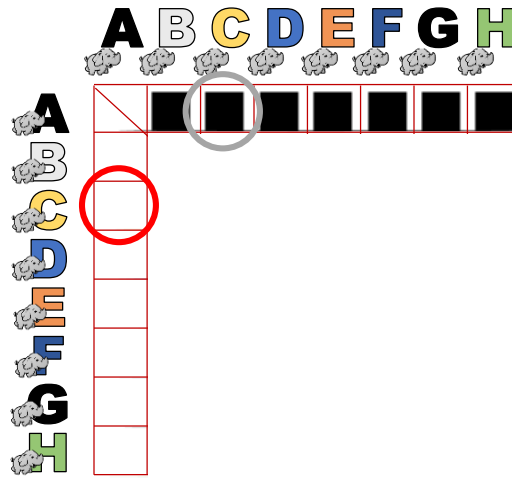
**DataTransferProtocol**

HDFS
DataNode

```
From readBlock In DataNode.DataTransferProtocol
  Join blockLocations In NameNode.GetBlockLocations
              On blockLocations -> readBlock
  Join cl In Client.DoRandomRead
              On cl -> blockLocations
  Where cl.host != readBlock.host
  GroupBy blockLocations.replicas, readBlock.host
  Select blockLocations.replicas, readBlock.host, COUNT
```
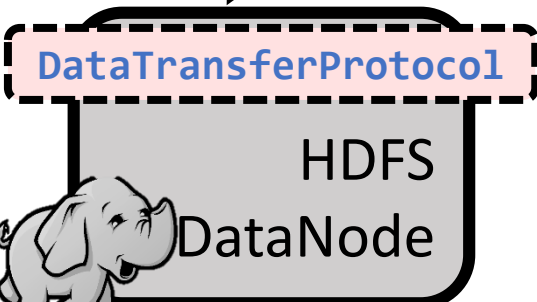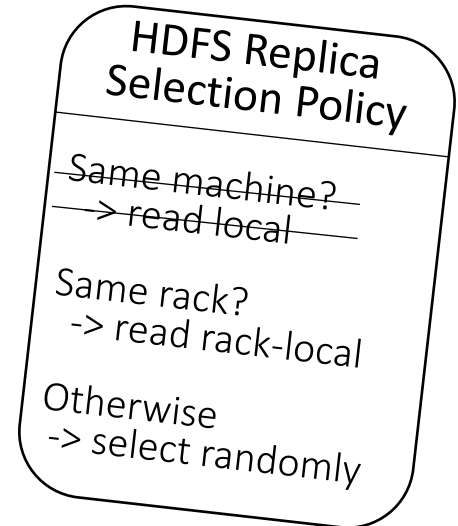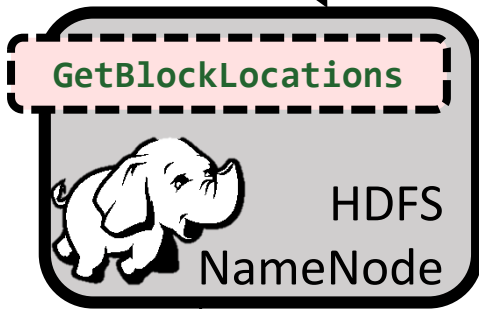
43

A B C D E F G H

DataNode Throughput [ops/s]

When both 🐘B and 🐘D host replicas,
Clients choose 🐘B this often: ___ (0%)

HDFS Replica Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

**Client**

**DoRandomRead**

**GetBlockLocations**

HDFS
NameNode

**DataTransferProtocol**

HDFS
DataNode

```
From readBlock In DataNode.DataTransferProtocol
    Join blockLocations In NameNode.GetBlockLocations
                    On blockLocations -> readBlock
    Join cl In Client.DoRandomRead
                    On cl -> blockLocations
    Where cl.host != readBlock.host
    GroupBy blockLocations.replicas, readBlock.host
    Select blockLocations.replicas, readBlock.host, COUNT
```
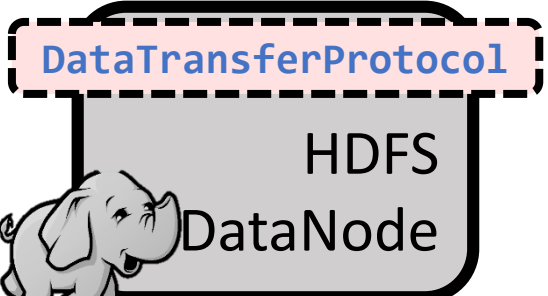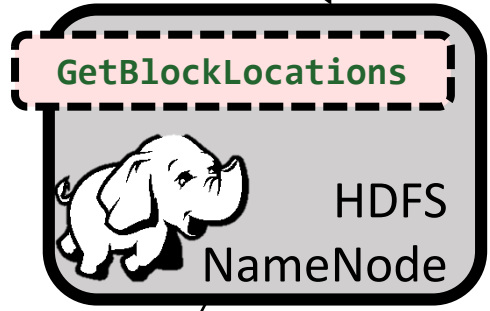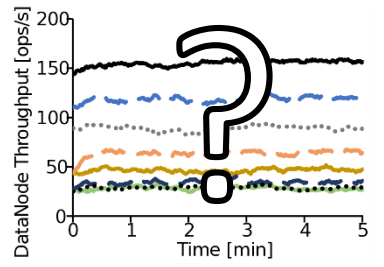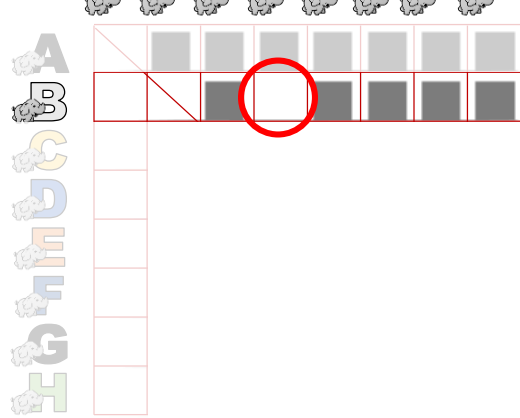
When both 🐘B and 🐘D host replicas,
Clients choose 🐘B this often: ___ (0%)

HDFS Replica
Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

Client

**DoRandomRead**

**GetBlockLocations**

HDFS
NameNode

**DataTransferProtocol**

HDFS
DataNode

```
From readBlock In DataNode.DataTransferProtocol
  Join blockLocations In NameNode.GetBlockLocations
                On blockLocations -> readBlock
  Join cl In Client.DoRandomRead
                On cl -> blockLocations
  Where cl.host != readBlock.host
  GroupBy blockLocations.replicas, readBlock.host
  Select blockLocations.replicas, readBlock.host, COUNT
```
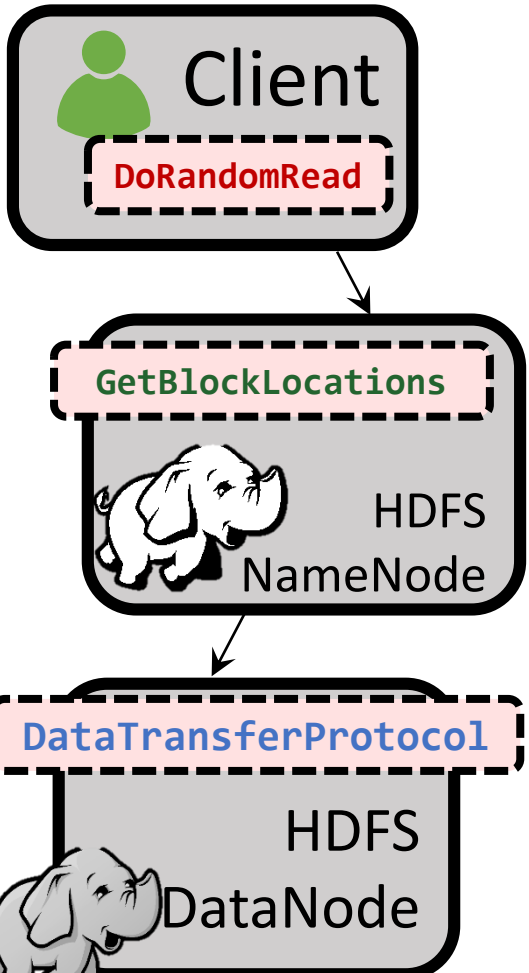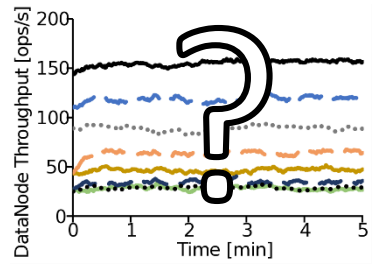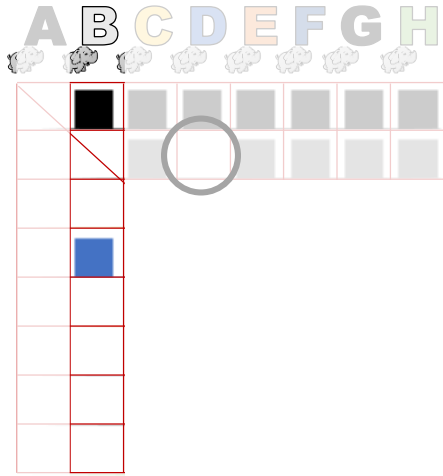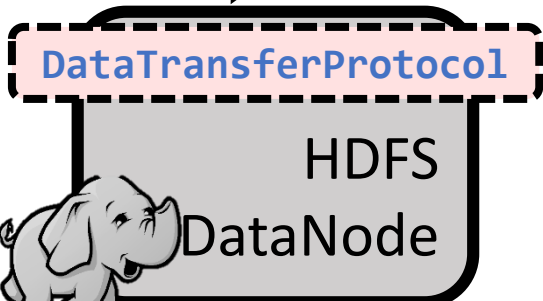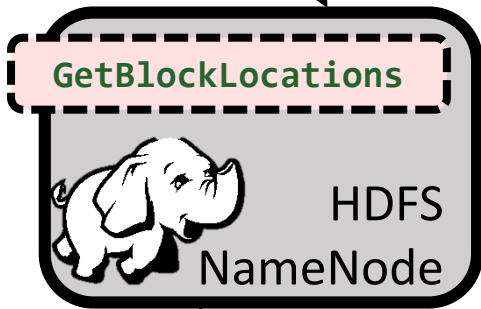
A B C D E F G H

When both B and D host replicas,
Clients choose B this often: ___ (0%)
Clients choose D this often: ■ (100%)

DataNode Throughput [ops/s]

200
150
100
50
0
0  1  2  3  4  5
Time [min]

**HDFS Replica Selection Policy**

~~Same machine? -> read local~~

Same rack? -> read rack-local

Otherwise -> select randomly

**Client**

**DoRandomRead**

**GetBlockLocations**

HDFS NameNode

**DataTransferProtocol**

HDFS DataNode

```
From readBlock In DataNode.DataTransferProtocol
   Join blockLocations In NameNode.GetBlockLocations
                  On blockLocations -> readBlock
   Join cl In Client.DoRandomRead
                  On cl -> blockLocations
   Where cl.host != readBlock.host
   GroupBy blockLocations.replicas, readBlock.host
   Select blockLocations.replicas, readBlock.host, COUNT
```
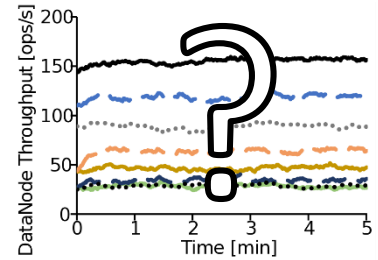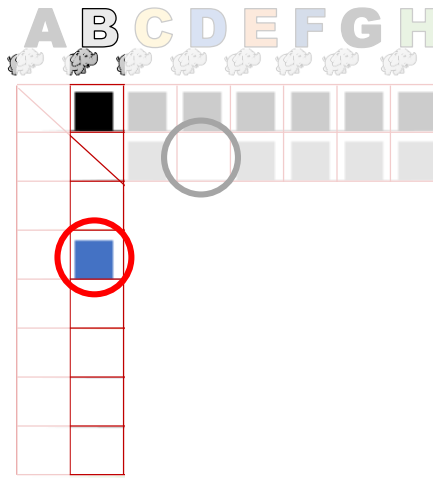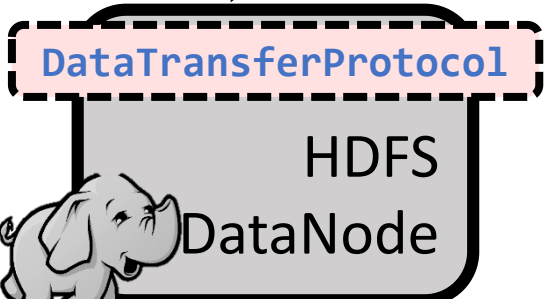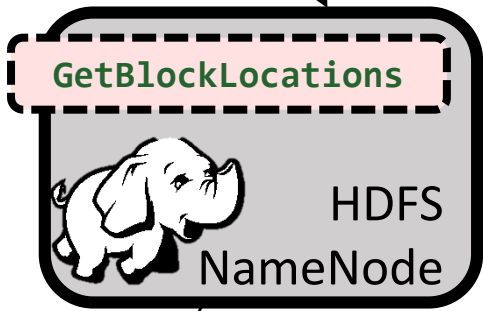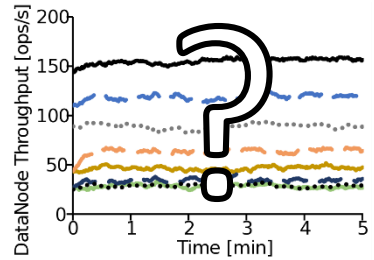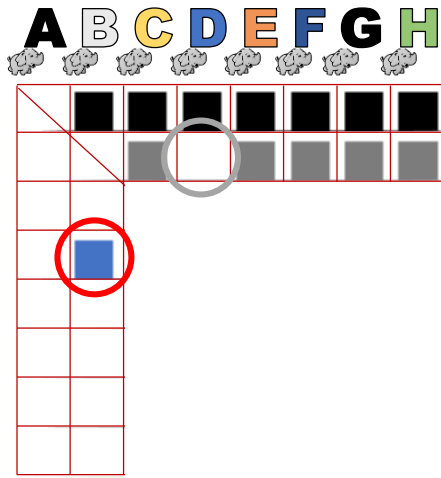
43

A B C D E F G H

When both B and D host replicas,
Clients choose B this often: ___ (0%)
Clients choose D this often: ■ (100%)

HDFS Replica Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

DataNode Throughput [ops/s]
200
150
100
50
0
0   1   2   3   4   5
Time [min]

Client

**DoRandomRead**

**GetBlockLocations**

HDFS NameNode

**DataTransferProtocol**

HDFS DataNode

```
From readBlock In DataNode.DataTransferProtocol
    Join blockLocations In NameNode.GetBlockLocations
                On blockLocations -> readBlock
    Join cl In Client.DoRandomRead
                On cl -> blockLocations
    Where cl.host != readBlock.host
    GroupBy blockLocations.replicas, readBlock.host
    Select blockLocations.replicas, readBlock.host, COUNT
```

A B C D E F G H

When both 🐘B and 🐘D host replicas,
Clients choose 🐘B this often: ___ (0%)
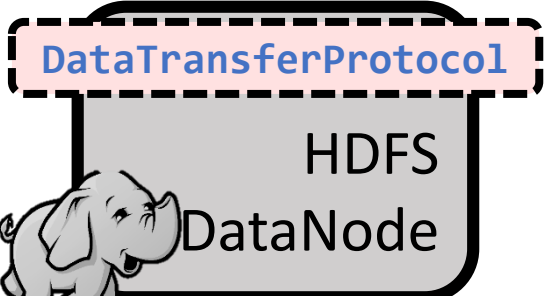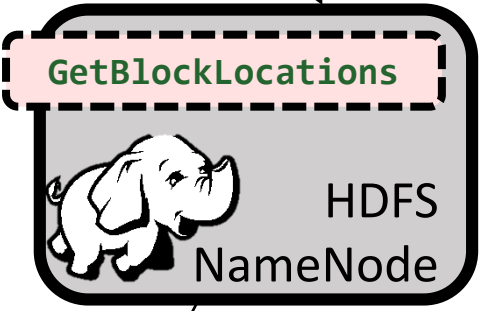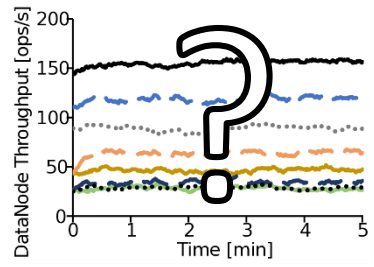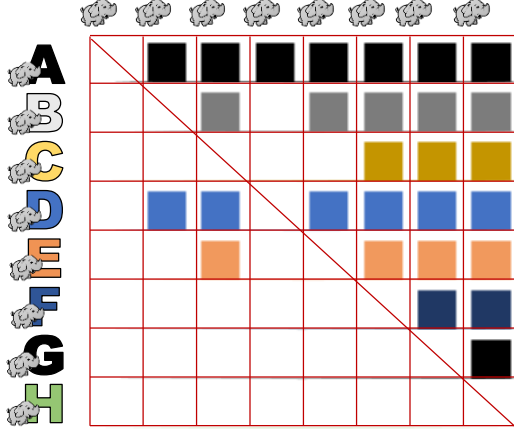Clients choose 🐘D this often: ■ (100%)

DataNode Throughput [ops/s]

Time [min]

**HDFS Replica Selection Policy**

~~Same machine? -> read local~~

Same rack? -> read rack-local

Otherwise -> select randomly

**Client**

**DoRandomRead**

**GetBlockLocations**

HDFS NameNode

**DataTransferProtocol**

HDFS DataNode

```
From readBlock In DataNode.DataTransferProtocol
    Join blockLocations In NameNode.GetBlockLocations
                On blockLocations -> readBlock
    Join cl In Client.DoRandomRead
                On cl -> blockLocations
    Where cl.host != readBlock.host
    GroupBy blockLocations.replicas, readBlock.host
    Select blockLocations.replicas, readBlock.host, COUNT
```
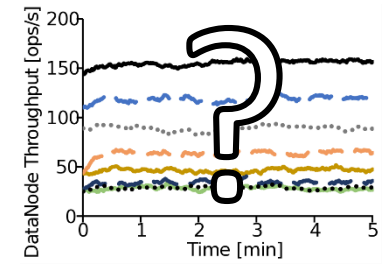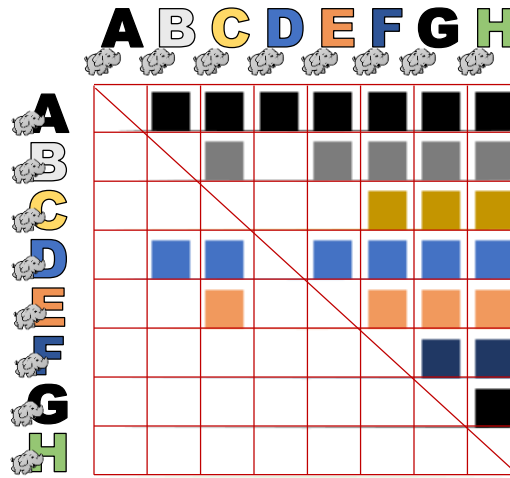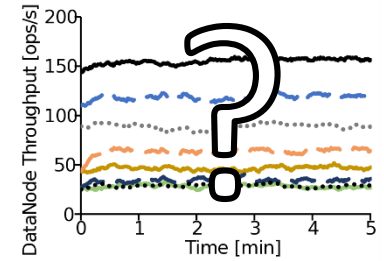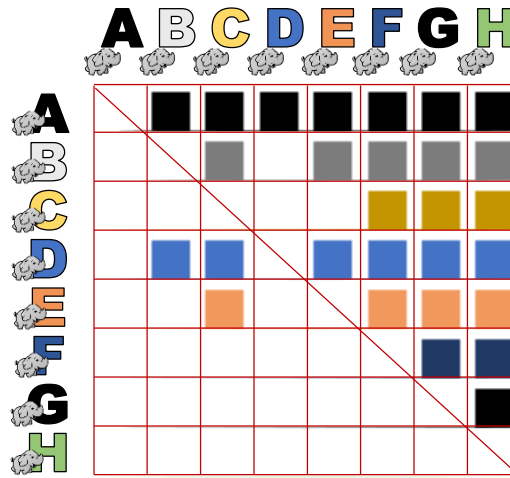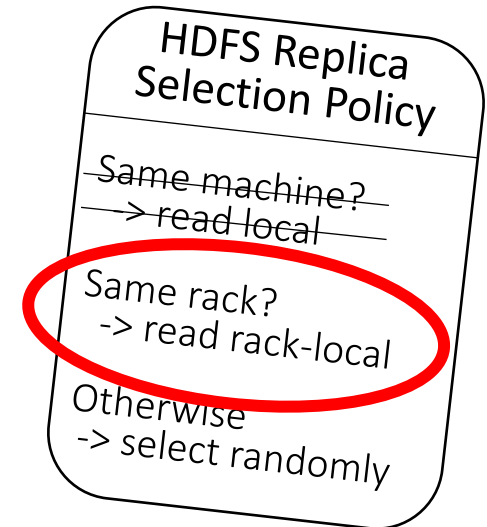
43

A B C D E F G H

DataNode Throughput [ops/s]

HDFS Replica
Selection Policy

~~Same machine?~~
~~-> read local~~

Same rack?
-> read rack-local

Otherwise
-> select randomly

- Lack of randomization skewed workload toward certain DNs

HDFS Replica
Selection Policy

Same machine?
-> read local

Same rack?
-> read rack-local

Otherwise
-> select randomly

- Lack of randomization skewed workload toward certain DNs

- **HDFS-6268**  Independently discovered. Fixed in HDFS 2.5

HDFS Replica Selection Policy

~~Same machine?~~
~~-> read local~~

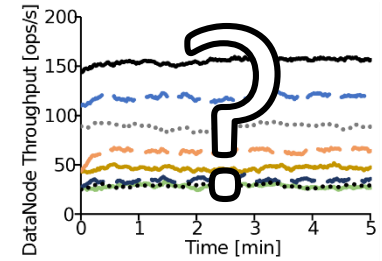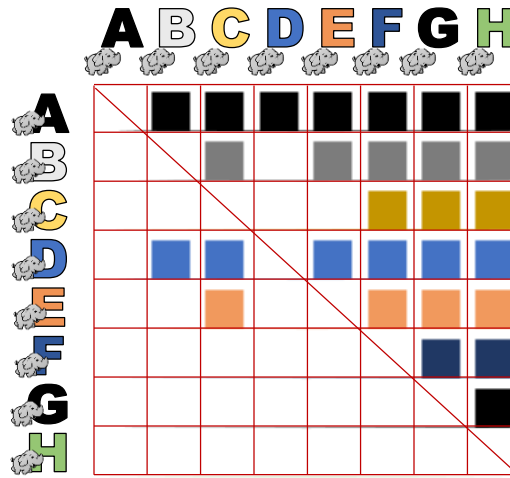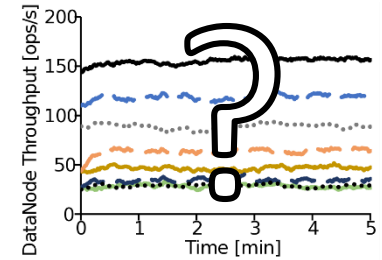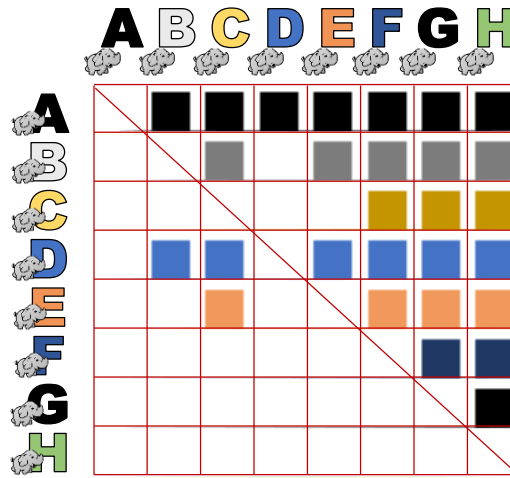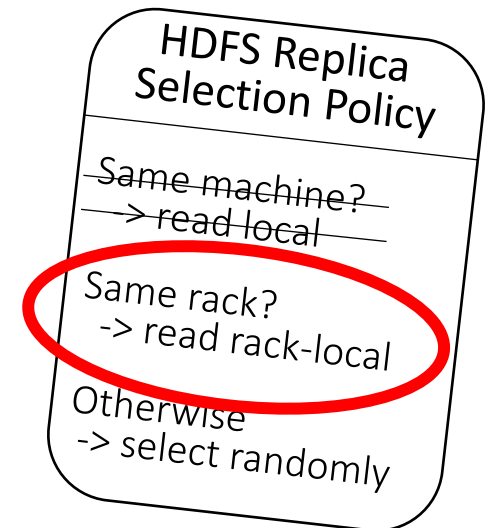Same rack?
-> read rack-local

Otherwise
-> select randomly

HDFS Replica Selection Policy

~~Same machine? -> read local~~

Same rack? -> read rack-local

Otherwise -> select randomly

- Lack of randomization skewed workload toward certain DNs

- **HDFS-6268** Independently discovered. Fixed in HDFS 2.5

- Seamlessly add correlations between multiple components
- Very specific, one-off metrics
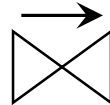- This experiment: 1.5% application-level overhead

# Pivot Tracing

Dynamic Causal Monitoring for Distributed Systems
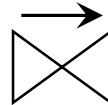
# Pivot Tracing

Dynamic Causal Monitoring for Distributed Systems

**Happened-Before Join**

# Pivot Tracing

Dynamic Causal Monitoring for Distributed Systems
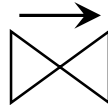
**Happened-Before Join**

Dynamic Instrumentation          Causal Tracing

# Pivot Tracing

Dynamic Causal Monitoring for Distributed Systems

**Happened-Before Join**

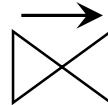Dynamic Instrumentation          Causal Tracing

Acceptable overheads for production (we think)

# Pivot Tracing

Dynamic Causal Monitoring for Distributed Systems

**Happened-Before Join**

Dynamic Instrumentation          Causal Tracing
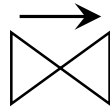
Acceptable overheads for production (we think)

Standing basic queries          Potential to dig deeper

# Pivot Tracing

## Dynamic Causal Monitoring for Distributed Systems

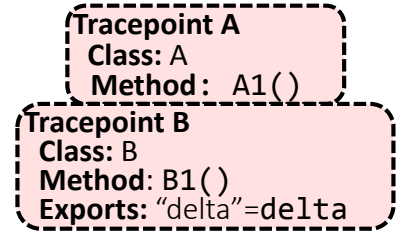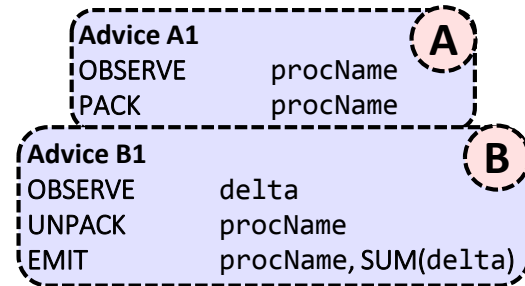**Happened-Before Join**

Dynamic Instrumentation

Causal Tracing

Acceptable overheads for production (we think)

Standing basic queries     Potential to dig deeper

**Tracepoint A**
**Class**: A
**Method**: A1()

**Tracepoint B**
**Class**: B
**Method**: B1()
**Exports**: "delta"=delta

```
From a In A
  Join b In B On a -> b
  GroupBy a.procName
  Select a.procName, SUM(b.delta)
```

**Advice A1**                      A
OBSERVE          procName
PACK             procName

**Advice B1**                      B
OBSERVE          delta
UNPACK           procName
EMIT             procName, SUM(delta)

**Jonathan Mace**     Ryan Roelke     Rodrigo Fonseca

BROWN
46