# Data Processing at the Speed of 100 Gbps using Apache Crail

Patrick Stuedi
IBM Research

# Apache Crail (crail.apache.org)



Apache Crail (Incubating) is a high-performance distributed data store designed for fast sharing of ephemeral data in distributed data processing workloads
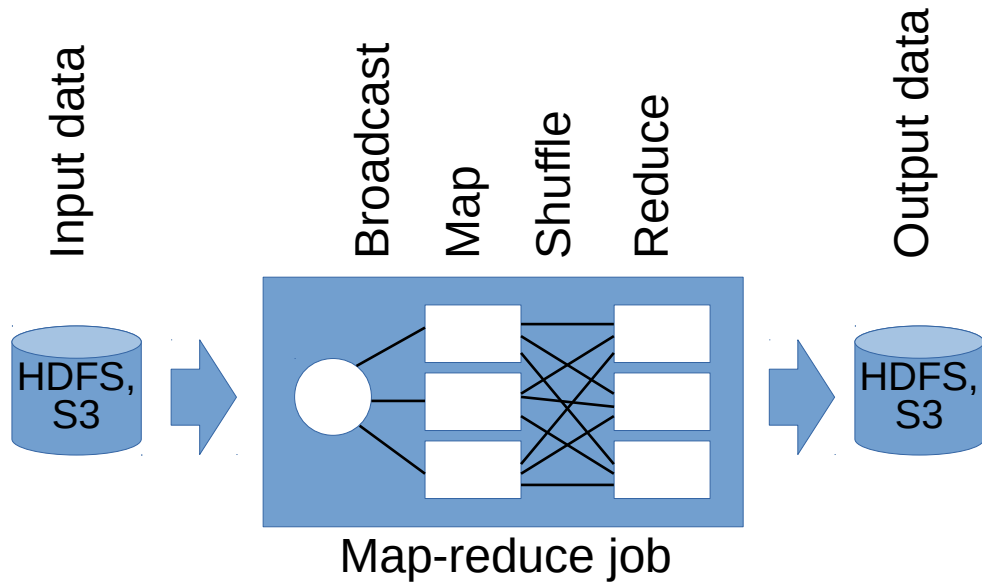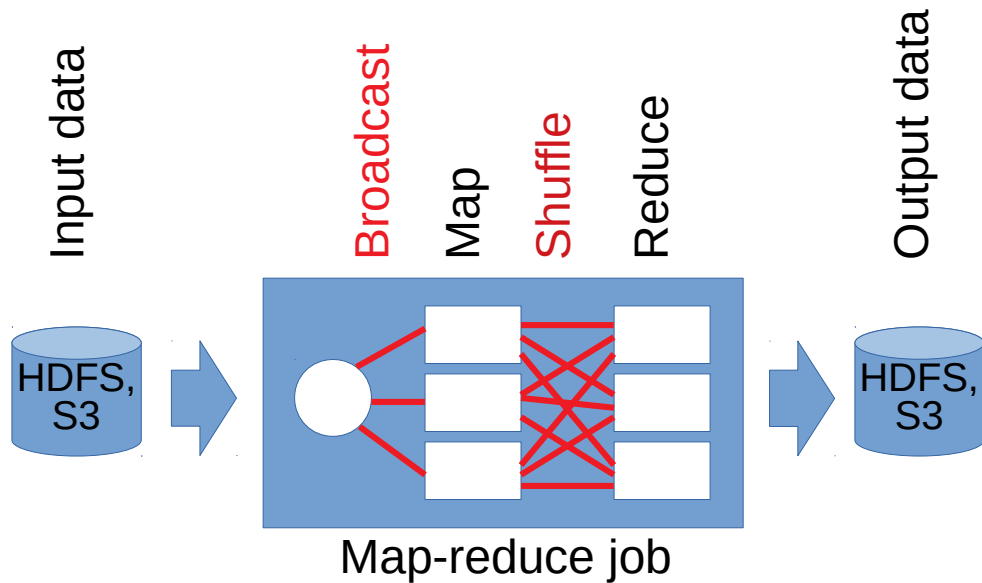
Download Now

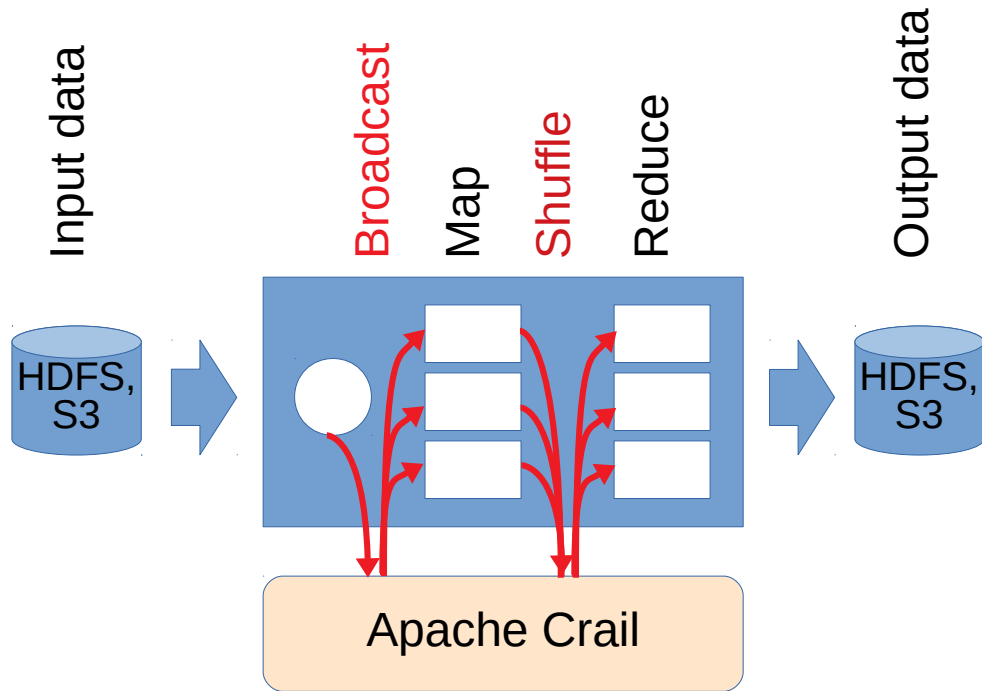Home   Overview   Downloads   Blog   Community   Documentation

# Apache Crail (crail.apache.org)

# Ephemeral Data

Input data

Broadcast

Map

Shuffle

Reduce

Output data

HDFS, S3

HDFS, S3

Map-reduce job

# Ephemeral Data

# Ephemeral Data

# Ephemeral Data

Input data

Broadcast

Map

Shuffle

Reduce

Intermediate data

HDFS, S3

HDFS, S3

HDFS, S3

Apache Crail

# Ephemeral Data

# Ephemeral Data

ML pre-processing
(map-reduce job)

normalized
images

ML training
(Tensorflow job)

Input data



HDFS, S3

HDFS, S3

HDFS, S3

Apache Crail

# Ephemeral Data

ML pre-processing
(map-reduce job)

normalized
images

ML training
(Tensorflow job)

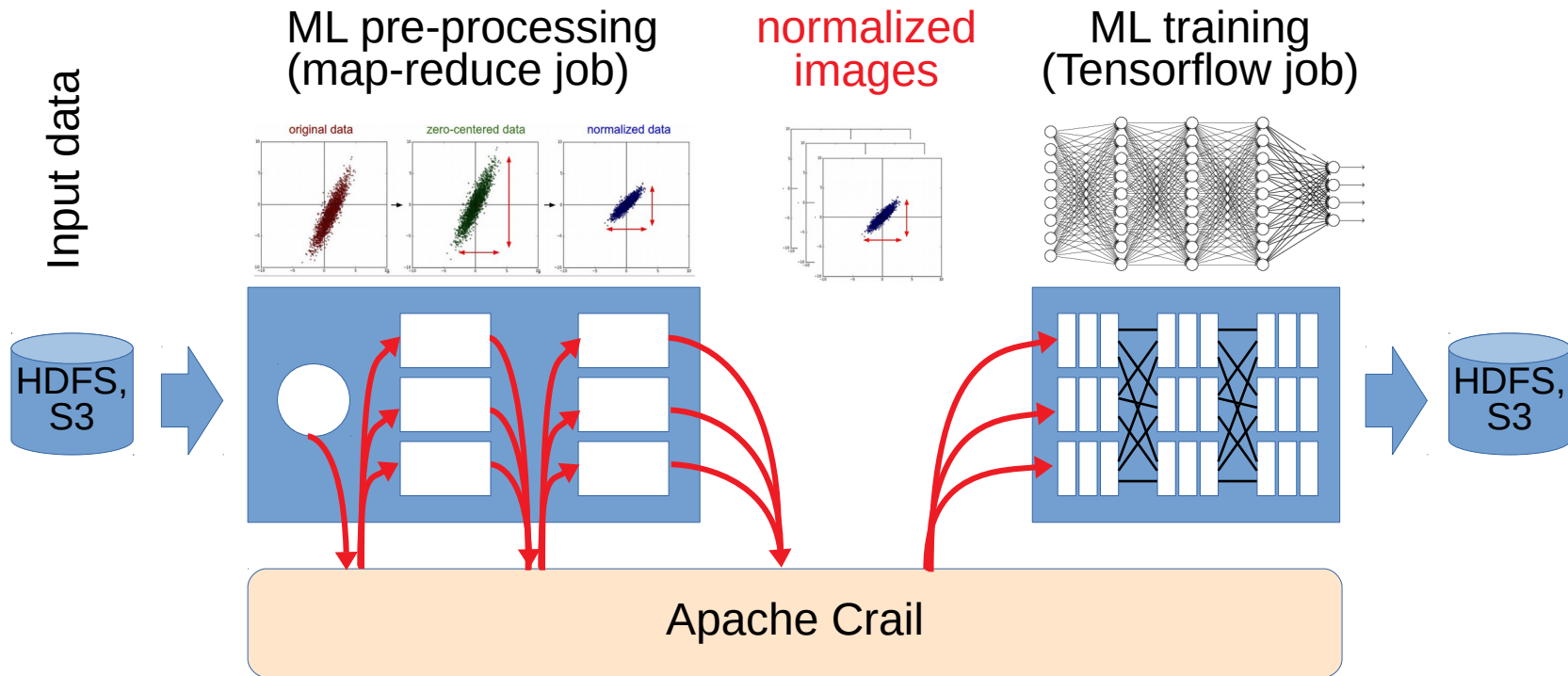Input data

original data     zero-centered data     normalized data

HDFS,
S3

HDFS,
S3

Apache Crail

# Why/when to use Crail



HDD

10 Gb/s Ethernet

# Why/when to use Crail

No
Crail
needed

HDD

10 Gb/s Ethernet

8 TB

100MB/s
10ms

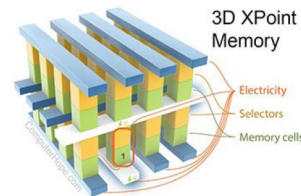10Gb/s
20us

# Why/when to use Crail



No Crail needed

100MB/s
10ms

10Gb/s
20us

HDD

Ultrastar
DC HC320

8 TB

10 Gb/s Ethernet

100x

10GB/s
10us

200Gb/s
1us

memory
data
RDMA
system network
latency

3D XPoint Memory

Electricity
Selectors
Memory cells

100 Gb/s Infiniband

nvm
EXPRESS®

Crail land

# Why/when to use Crail



No Crail needed

HDD
Ultrastar DC HC320
8TB

10 Gb/s Ethernet

100MB/s
10ms

10Gb/s
20us

100x

memory
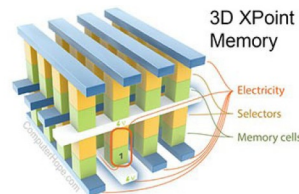data
RDMA
system
network

3D XPoint Memory
Electricity
Selectors
Memory cells

100 Gb/s Infiniband

nvm EXPRESS®

Crail land

10GB/s
10us

200Gb/s
1us

Throughput (Gbit/s)
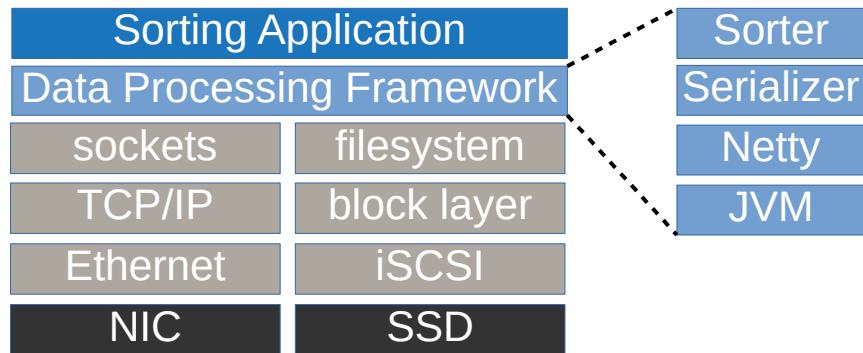
100
80
60
40
20
0

88.3s

hardware limit

Spark/Crail
Spark/Vanilla

527.6s

Elapsed time

Terasort
12.8 TB data
128 nodes

# Performance Challenge

| | 1 Gbps | HDD | 100 Gbps | Flash |
|---|---|---|---|---|
| Bandwidth | 117 MB/s | 140 MB/s | 12.5 GB/s | 3.1 GB/s |
| cycle/unit | 38,400 | 10,957 | 360 | 495 |

Sorting Application

Data Processing Framework

sockets | filesystem

TCP/IP | block layer

Ethernet | iSCSI

NIC | SSD

Sorter

Serializer

Netty

JVM

# Performance Challenge

| | 1 Gbps | HDD | 100 Gbps | Flash |
|---|---|---|---|---|
| Bandwidth | 117 MB/s | 140 MB/s | 12.5 GB/s | 3.1 GB/s |
| cycle/unit | 38,400 | 10,957 | 360 | 495 |

| Sorting Application |
|---|
| Data Processing Framework |

| sockets | filesystem |
|---|---|
| TCP/IP | block layer |
| Ethernet | iSCSI |
| NIC | SSD |

| Sorter |
|---|
| Serializer |
| Netty |
| JVM |

Process chunk
In reduce task



Fetch chunk
Over the network

HotNets'16

# Performance Challenge

| | 1 Gbps | HDD | 100 Gbps | Flash |
|---|---|---|---|---|
| Bandwidth | 117 MB/s | 140 MB/s | 12.5 GB/s | 3.1 GB/s |
| cycle/unit | 38,400 | 10,957 | 360 | 495 |



HotNets'16

# Performance Challenge

| | 1 Gbps | HDD | 100 Gbps | Flash |
|---|---|---|---|---|
| Bandwidth | 117 MB/s | 140 MB/s | 12.5 GB/s | 3.1 GB/s |
| cycle/unit | 38,400 | 10,957 | 360 | 495 |

software overhead are spread over the entire stack

| Sorting Application |  |
|---|---|
| Data Processing Framework |  |
| sockets | filesystem |
| TCP/IP | block layer |
| Ethernet | iSCSI |
| NIC | SSD |

| Sorter |
| Serializer |
| Netty |
| JVM |



- Misc.
- Iterator
- Ser
- Sorting
- IO
- JVM
- Linux

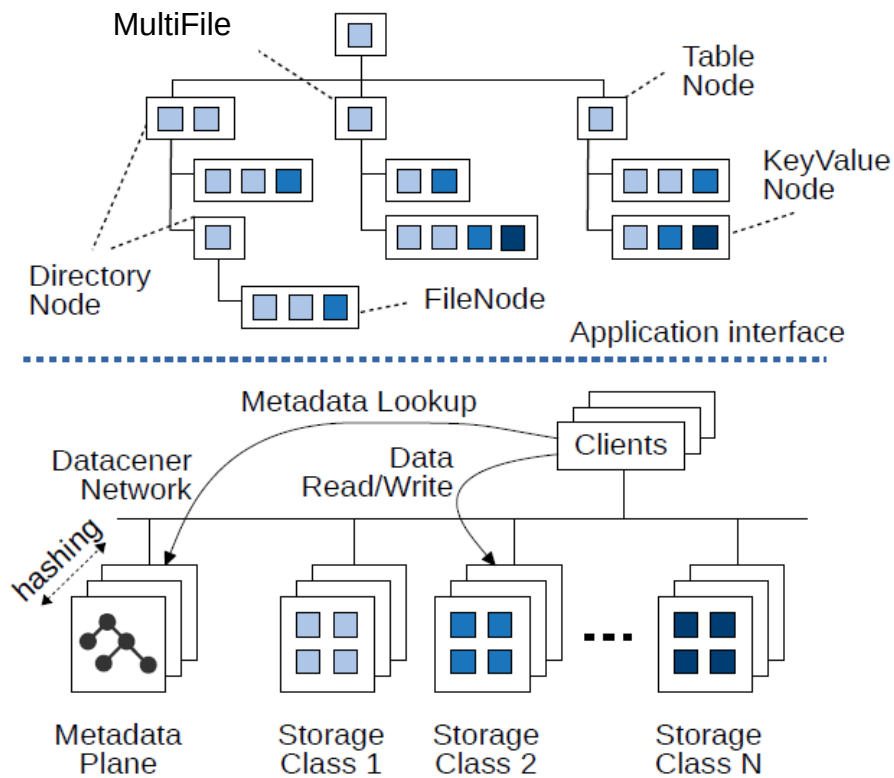1 Gbps: 52%, 48%
10 Gbps: 92%, 8%
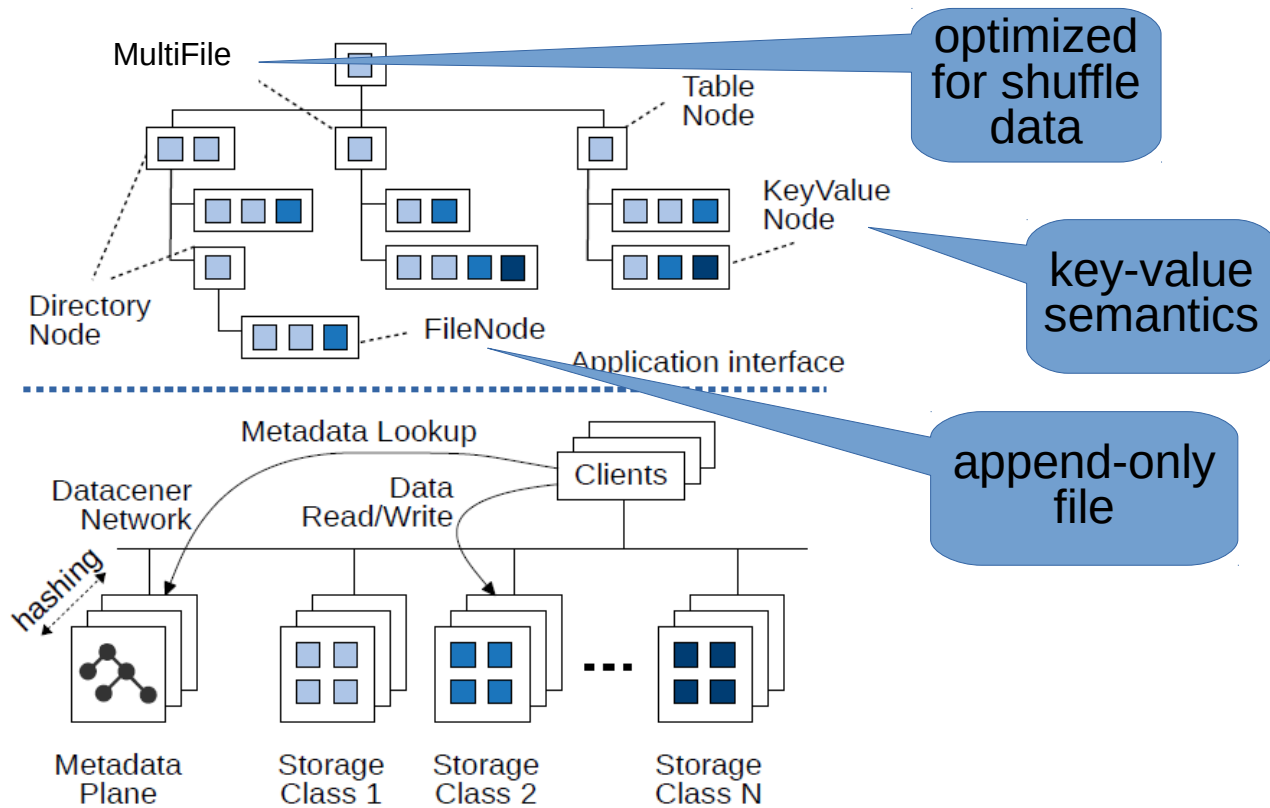40 Gbps: 97%

HotNets'16
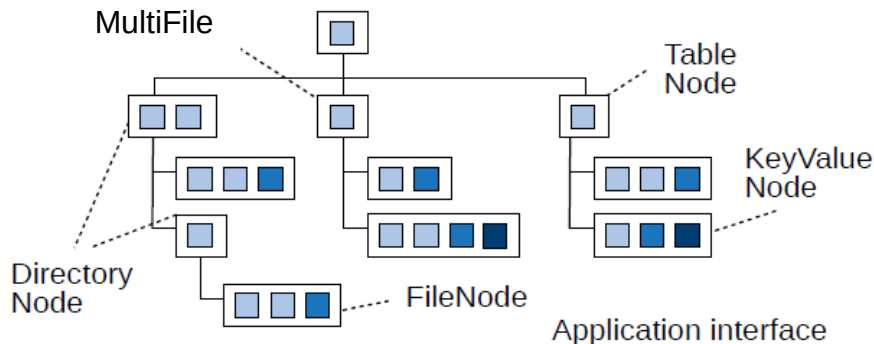
# Crail Overview

# Crail Overview

# Crail Architecture & API
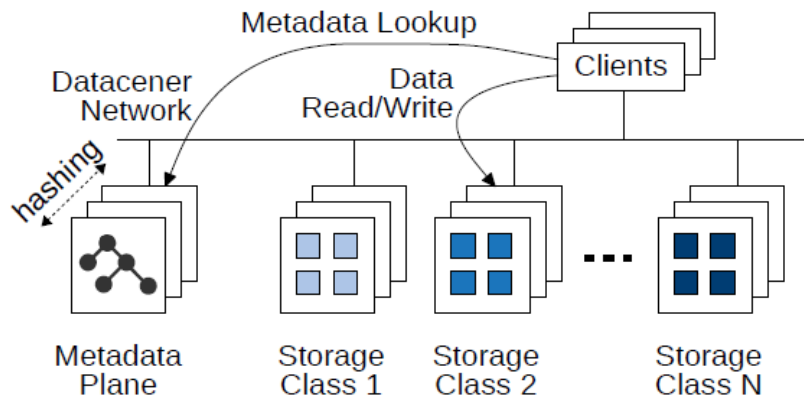
# Crail Architecture & API
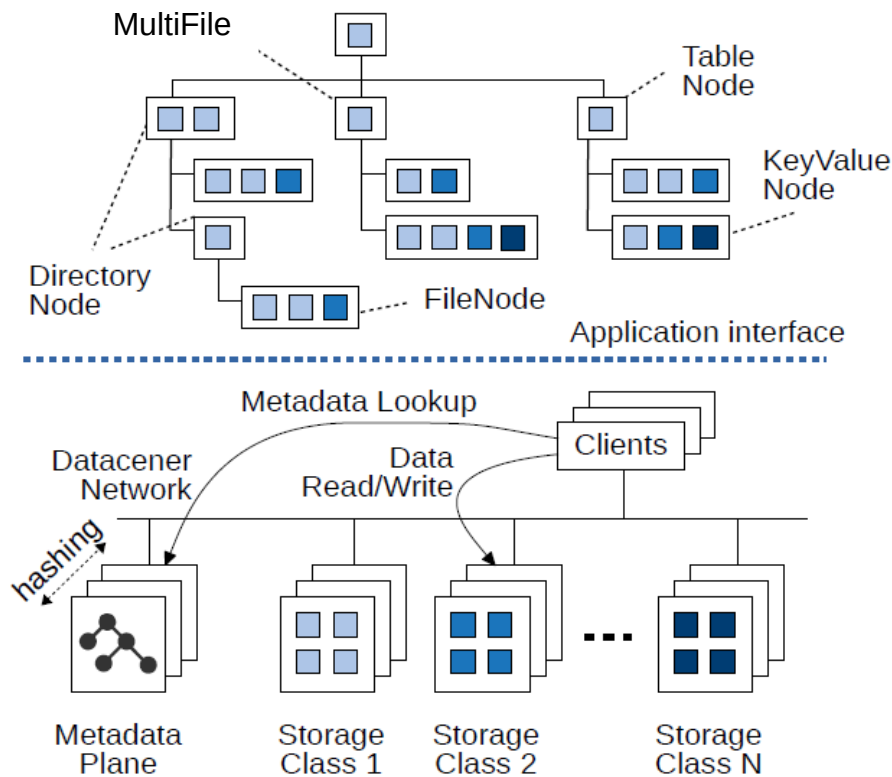
# Crail Architecture & API

Java:

```
CrailStore crail = CrailStore.newInstance();
Future<Node> fut = crail.create("/a.dat", CrailType.File);
//...do work
CrailFile file = fut.get().asFile();
CrailOutputStream stream = file.getDirectOutputStream();
ByteBuffer buffer = crail.allocateBuffer();
Future<CrailResult> ret = stream.write(buf);
//...do work
ret.get();
```

C++:

```
CrailStore crail;
auto fut = crail.Create<CrailFile>("/tmp.dat");
//..do work
CrailFile file = fut.get();
CrailOutputStream stream = file.outputstream();
shared_ptr<ByteBuffer> buf = make_shared<ByteBuffer>(len);
Future<int> ret = stream.Write(buf);
//..do work
ret.get();
```
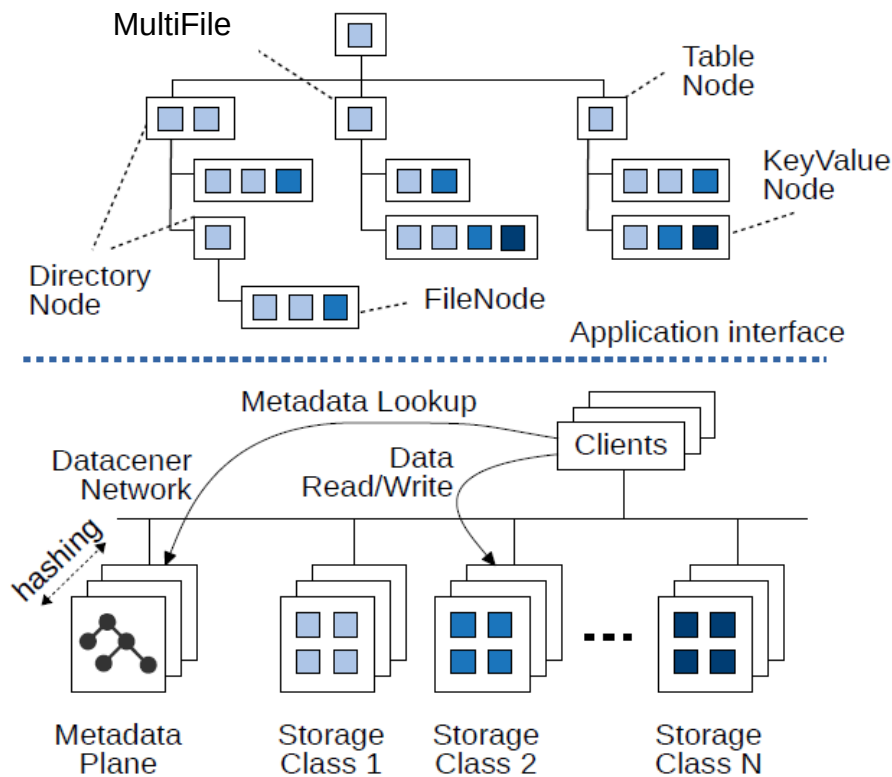
# Crail Architecture & API

```java
CrailStore crail = CrailStore.newInstance();
Future<Node> fut = crail.create("/a.dat", CrailType.File);
//...do work
CrailFile file = fut.get().asFile();
CrailOutputStream stream = file.getDirectOutputStream();
ByteBuffer buffer = crail.allocateBuffer();
Future<CrailResult> ret = stream.write(buf);
//...do work
ret.get();
```

Node type

C++:

```cpp
CrailStore crail;
auto fut = crail.Create<CrailFile>("/tmp.dat");
//..do work
CrailFile file = fut.get();
CrailOutputStream stream = file.outputstream();
shared_ptr<ByteBuffer> buf = make_shared<ByteBuffer>(len);
Future<int> ret = stream.Write(buf);
//..do work
ret.get();
```

# Crail Architecture & API

```
CrailStore crail = CrailStore.newInstance();
Future<Node> fut = crail.create("/a.dat", CrailType.File);
//...do work
CrailFile file = fut.get().asFile();
CrailOutputStream stream = file.getDirectOutputStream();
ByteBuffer buffer = crail.allocateBuffer();
Future<CrailResult> ret = stream.write(buf);
//...do work
ret.get();
```
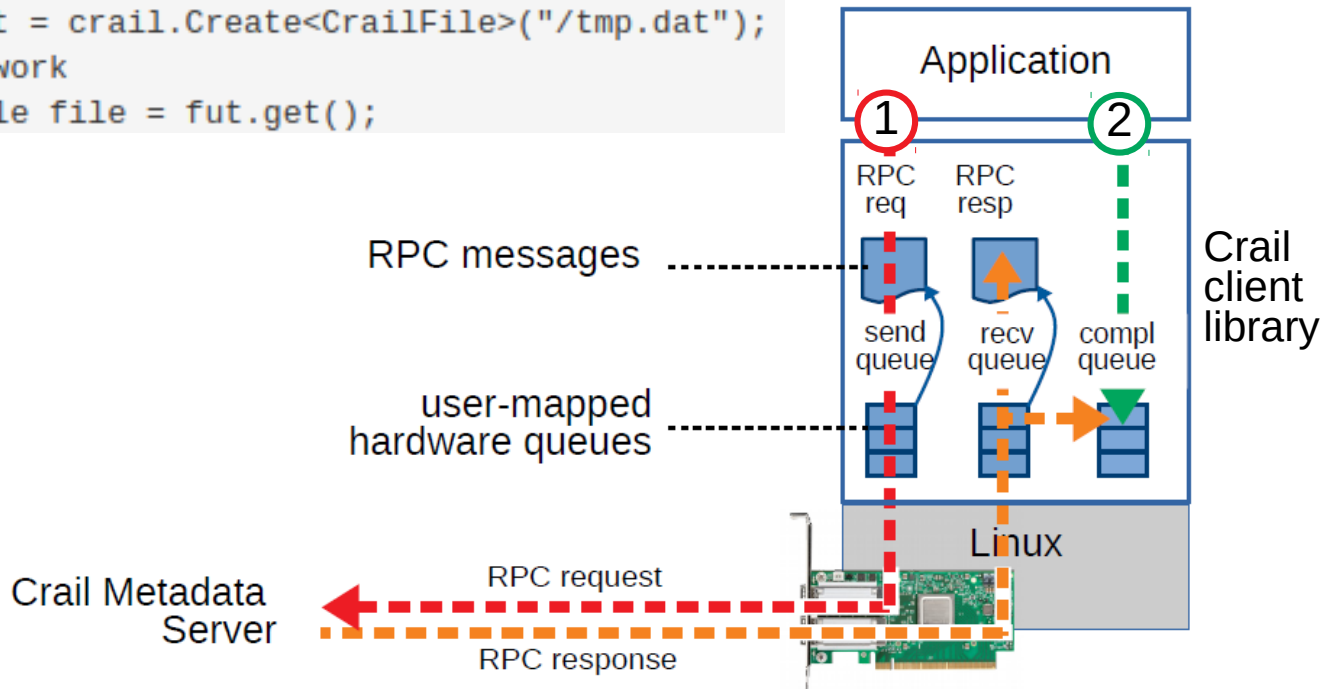
non-blocking & asynchronous

C++:

```
CrailStore crail;
auto fut = crail.Create<CrailFile>("/tmp.dat");
//..do work
CrailFile file = fut.get();
CrailOutputStream stream = file.outputstream();
shared_ptr<ByteBuffer> buf = make_shared<ByteBuffer>(len);
Future<int> ret = stream.Write(buf);
//..do work
ret.get();
```

# Where does the performance come from?
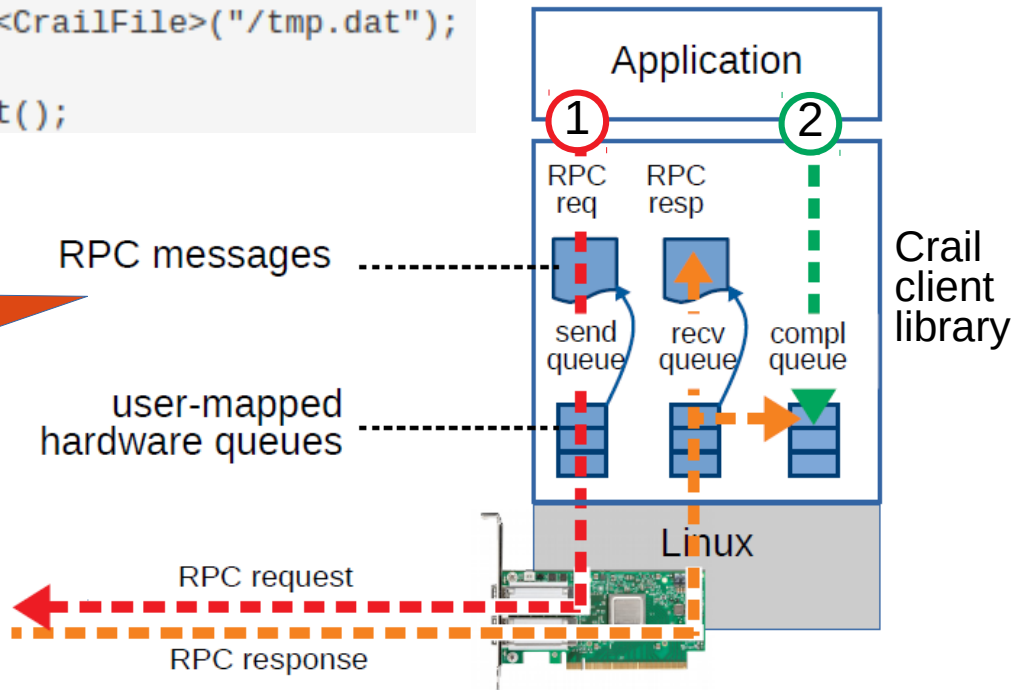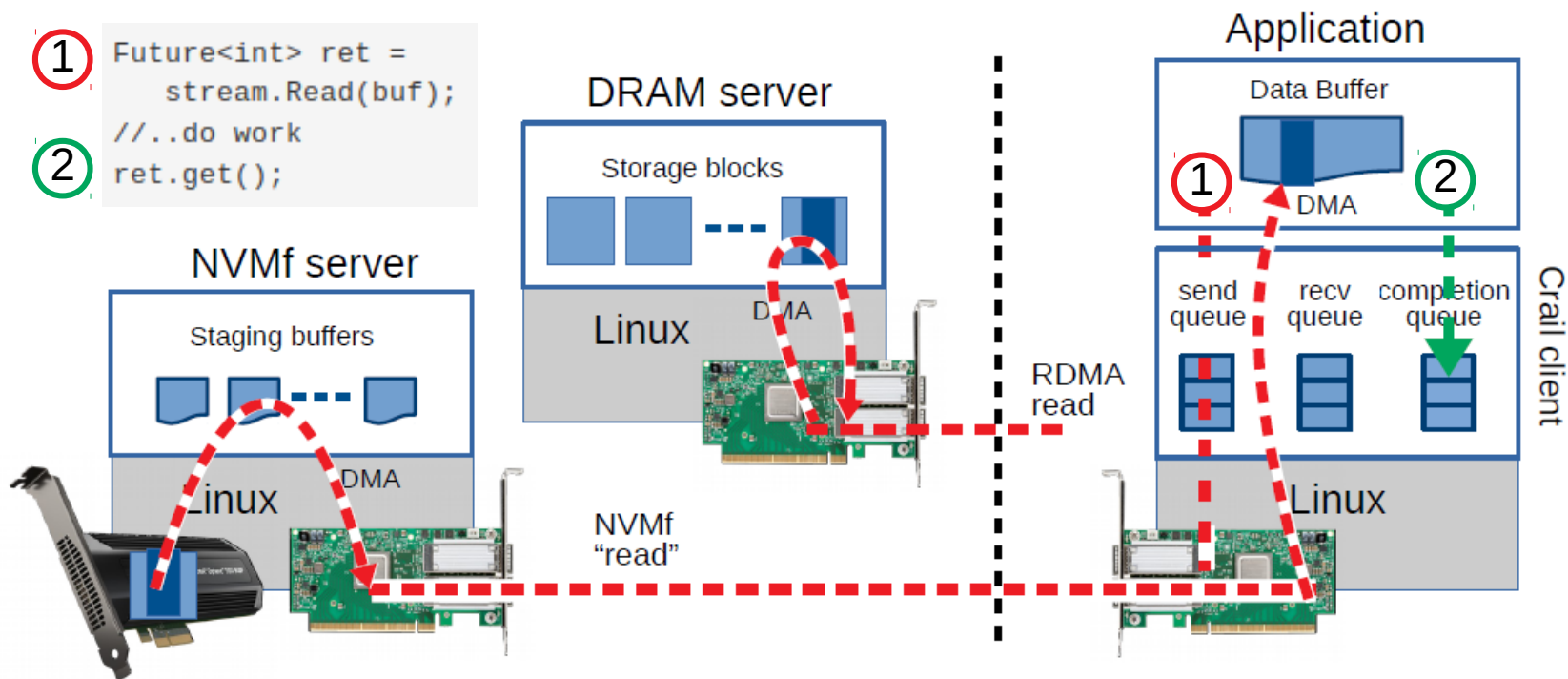
# User-Level I/O: Metadata

# User-Level I/O: Metadata

```
① auto fut = crail.Create<CrailFile>("/tmp.dat");
  //..do work
② CrailFile file = fut.get();
```
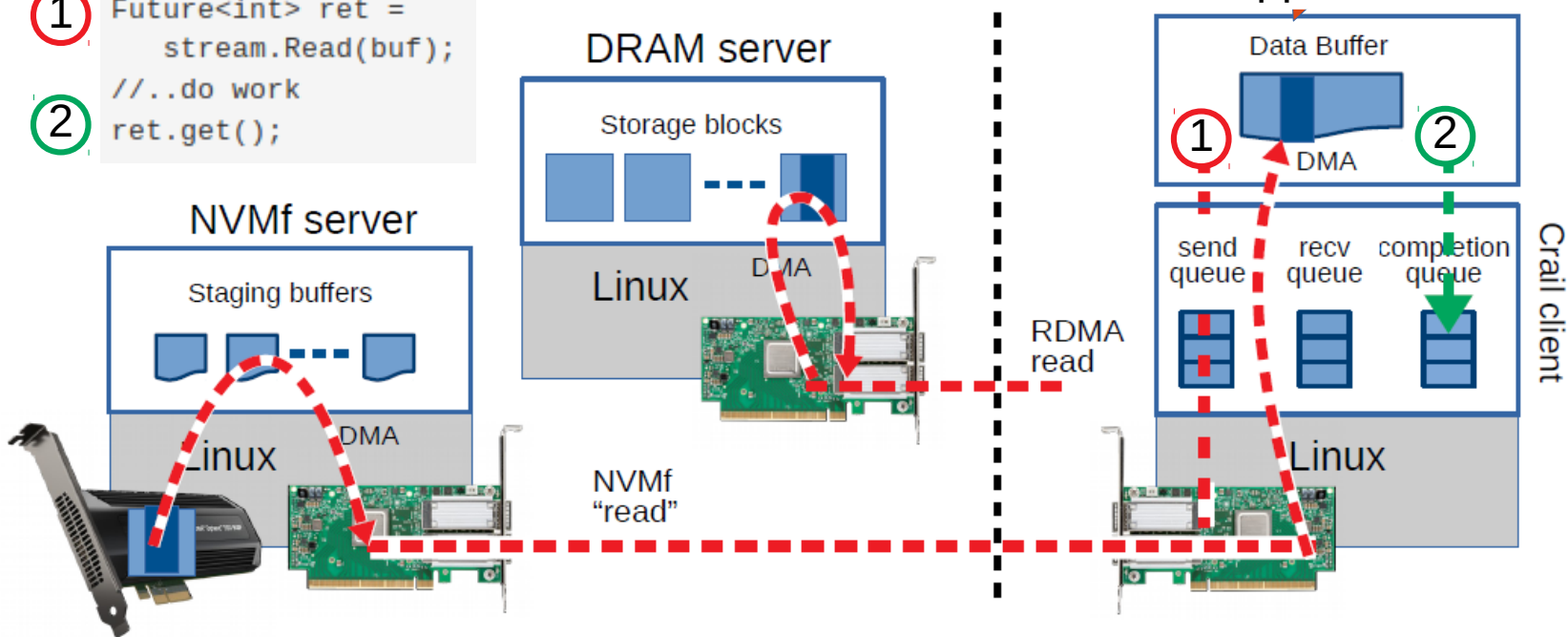
No threads
No context switches

Application

① ②

Crail client library

RPC req    RPC resp

RPC messages

send queue    recv queue    compl queue

user-mapped hardware queues

Linux

Crail Metadata Server

RPC request

RPC response

# User-Level I/O: Data

# User-Level I/O: Data

zero-copy, transfer only data that is requested
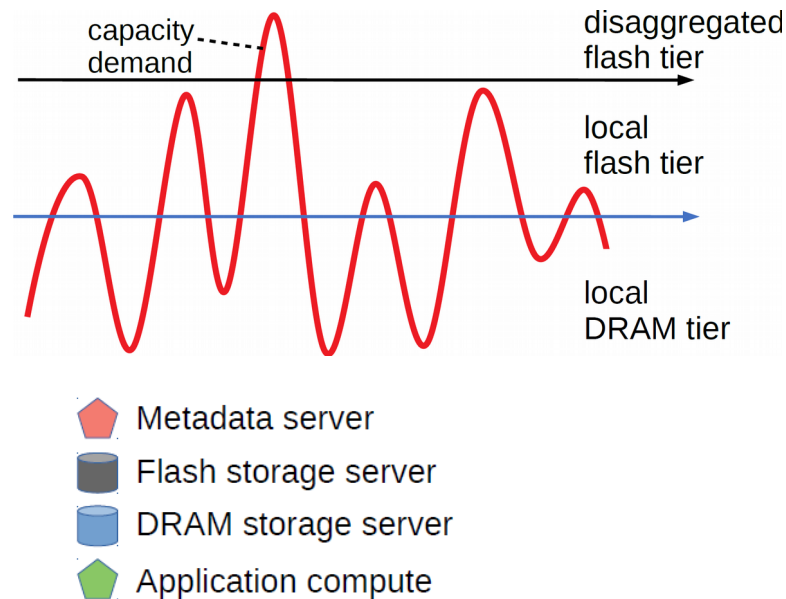
# Crail Deployment Modes



compute/storage
co-located

storage
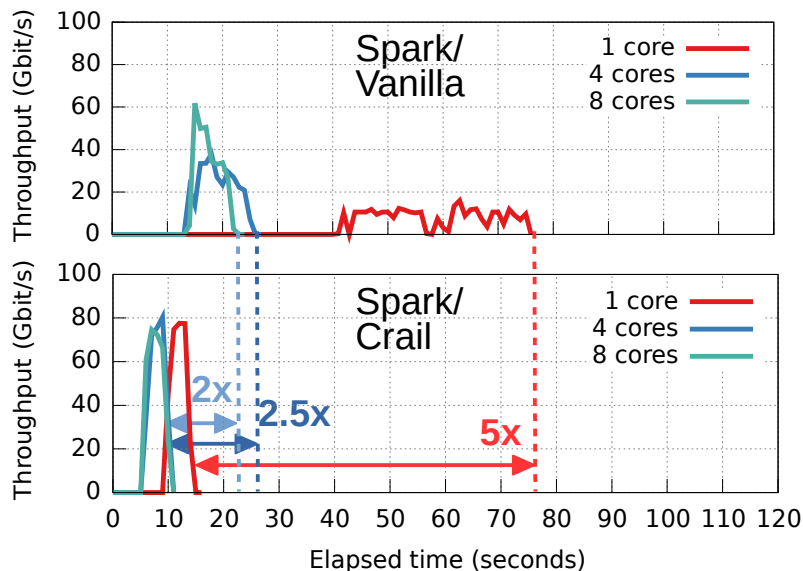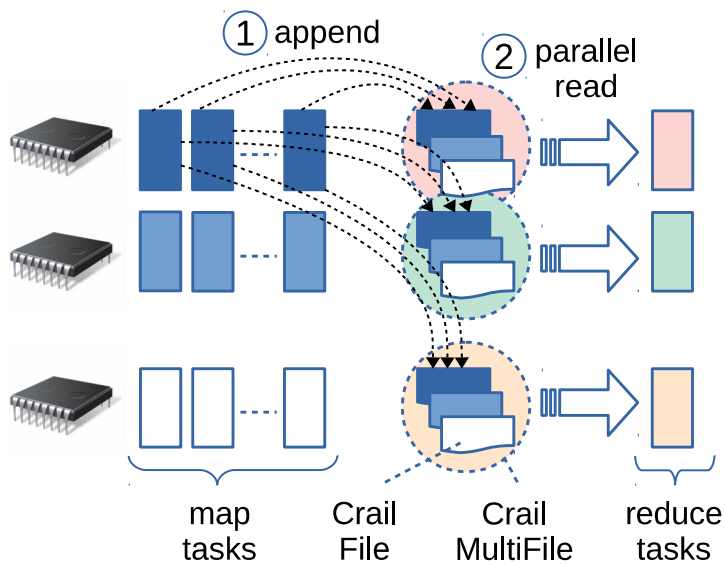disaggregation

flash storage
disaggregation

Metadata server
Flash storage server
DRAM storage server
Application compute

# YCSB KeyValue Workload
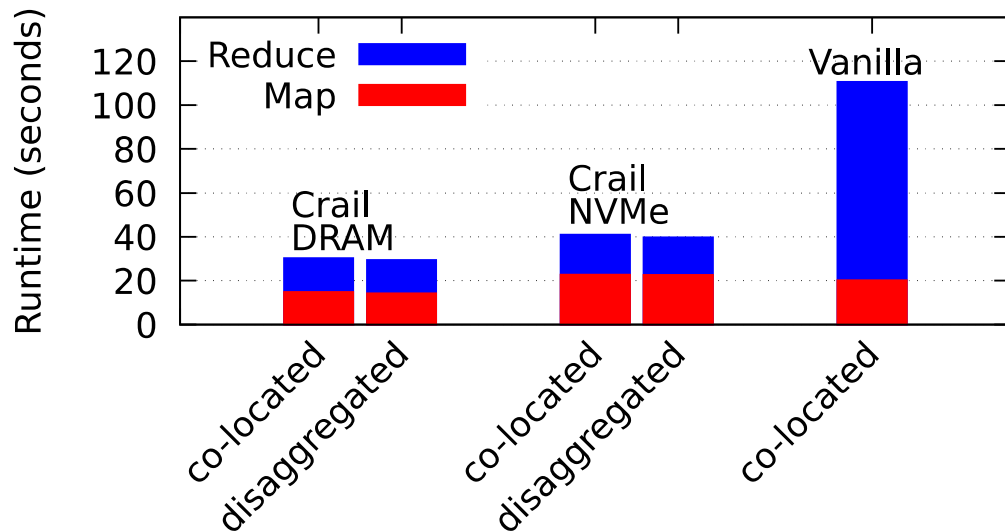


GET
Value size:
1KB

GET
Value size:
100KB

Crail offers Get latencies of ~12us and 30us for DRAM and NVM for 100 byte KV pairs
Crail offers Get latencies of ~30us and 40us for DRAM and NVM for 1000 byte KV pairs
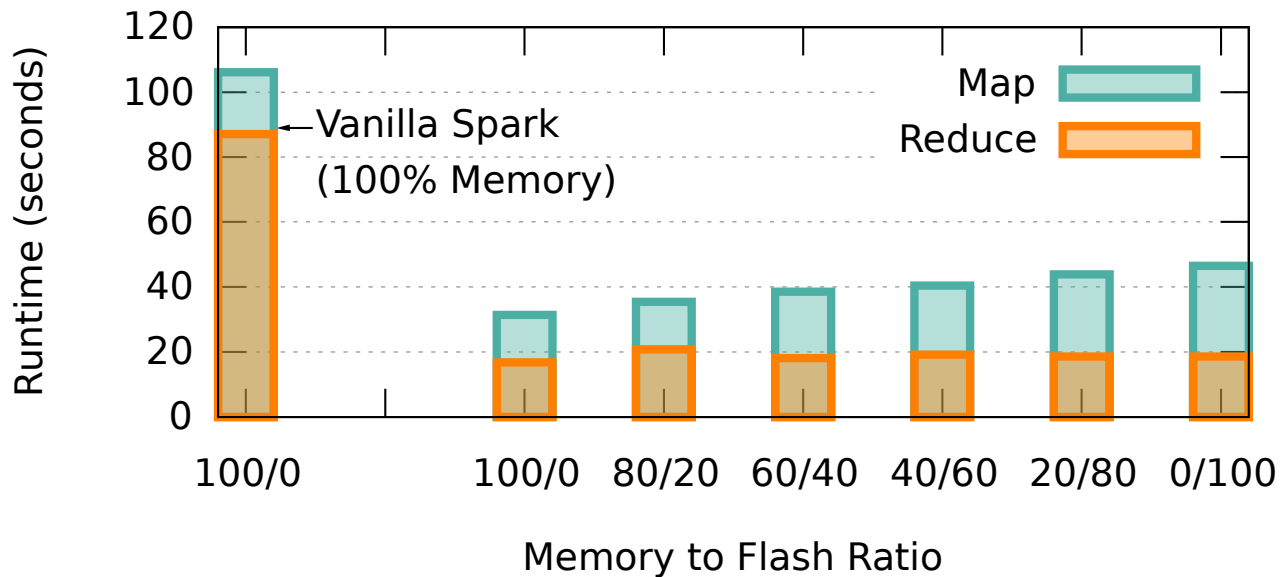
# Spark GroupBy (80M keys, 4K)



Spark shuffling via Crail on a single core is 2x faster
than vanilla Spark on 8 cores per executor (8 executors)

# DRAM & Flash Disaggregation



Crail enables disaggregation of temporary data at no cost

# DRAM/Flash Tiering



Using flash only increases the sorting time by around 48%

# Conclusions

- Apache Crail: Fast distributed "tmp"
  - User-level I/O
  - Storage disaggregation
  - Memory/flash convergence

- Applications
  - Intra-job scratch space (shuffle, broadcast, etc.)
  - Multi-job pipelines

- Coming soon
  - Native Crail (C++)
  - Tensorflow-Crail