# UNORTHODOX PATHS TO HIGH PERFORMANCE

@ALEXRAS

TRIFACTA

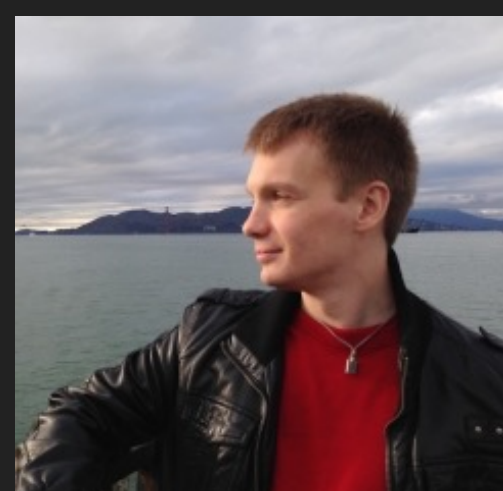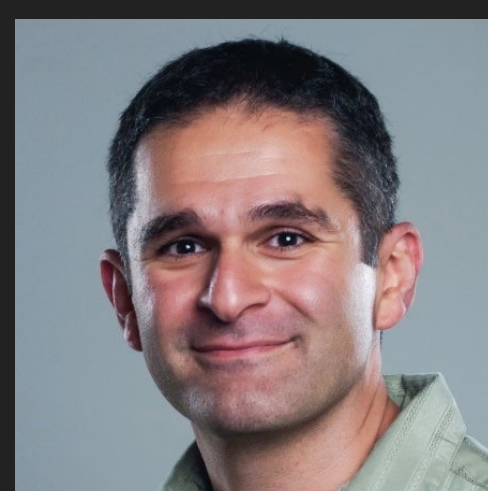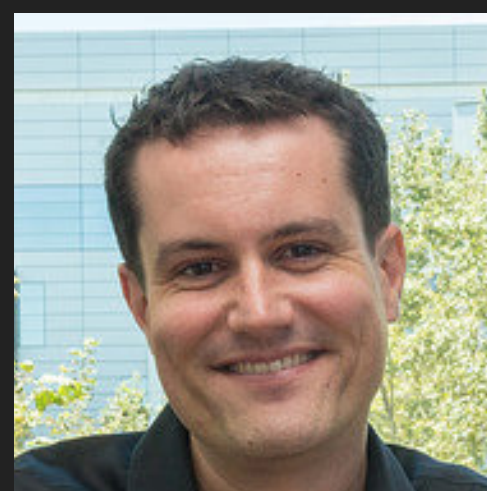# TRITONSORT (NSDI 2011)
Sort Really Fast

# THEMIS (SOCC 2012)
MapReduce Really Fast

# THIS TALK:
# HOW WE DID IT

# 1. SOFTWARE–HARDWARE CO-DESIGN

# 2. BUILDING FOR EXPERIMENTATION

# 3. CAREFULLY MANAGING MEMORY

MOTIVATION

[SORT] IS AN EXCELLENT TEST OF THE INPUT-OUTPUT ARCHITECTURE OF A COMPUTER AND ITS OPERATING SYSTEM.

"A measure of transaction processing power"
Datamation 1985

# THE SORT BENCHMARK

- SORTING K/V PAIRS (RECORDS)
- MANY CATEGORIES, VARIANTS
- TODAY: GRAYSORT (100TB)

2009: YAHOO! SORTS 100TB IN 173 MINUTES ON 3452 HADOOP NODES

**578** GB PER MINUTE

**9.6** GB PER SECOND

**3452** NODES

**2.79** MBPS PER NODE

THE GAP IN THEORETICAL PER NODE PERFORMANCE AND WHAT IS ACTUALLY ACHIEVED HAS BECOME GLARINGLY LARGE.

Anderson and Tucek
"Efficiency Matters!" SIGOPS  OSR 2010

# HOW CAN WE DO BETTER?

# HARDWARE & SOFTWARE CO-DESIGNED FOR WORKLOAD

UNDERSTANDING THE PROBLEM

# HOW TO MAXIMIZE PER-NODE SPEED?

~~CPU~~

~~RAM~~

NETWORK 10Gbps

DISK ~15 disks/NIC

▶ **8 CORES**  ▶ **24 GB RAM**

▶ **16 DISKS**  ▶ **10 GBPS NIC**

# WHAT ARE THE EXPENSIVE OPERATIONS?

🐌 **WRITING**

💀 **SEEKING**

# KEEP OFF THE DISK

## TWO READS + TWO WRITES PER RECORD

# SEEK INFREQUENTLY

## BIG READS
## BIG WRITES

# WHICH DISTRIBUTED SORTING ALGORITHM?

# MERGESORT

# DISTRIBUTION SORT

- **BIG WRITES** IN FIRST PASS
- **SEQUENTIAL I/O** IN SECOND PASS

NOW, TO BUILDING!

IT WON'T JUST MAGICALLY BE FAST

EXPERIMENTATION!

# FLEXIBLE
# MODULAR
# GRAPHS

STAGE

# DIFFERENT SORTS?

| READ | → | RADIX SORT | → | WRITE |

| READ | → | QUICKSORT | → | WRITE |

| READ | → | TIMSORT | → | WRITE |

MEASUREMENT

HOW LARGE ARE WRITES?

WHERE IS THE BOTTLENECK?

ARE STAGES BLOCKED? IDLE?

TONS AND TONS OF LOGS

ONE FULL CORD

# AGGREGATES



/mnt/disks/cciss/c1d0p1

/mnt/disks/cciss/c0d2p1

# TIME-SERIES

# RUNTIME INFO

**writer**

| Key | Worker ID | | | | | | | | Total |
|---|---|---|---|---|---|---|---|---|---|
| | **0** | **1** | **2** | **3** | **4** | **5** | **6** | **7** | |
| bytes_consumed | 7.88 GiB | 7.89 GiB | 7.89 GiB | 7.90 GiB | 7.87 GiB | 7.88 GiB | 7.89 GiB | 7.89 GiB | 63.09 GiB |
| bytes_produced | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B |
| bytes_written | 7.88 GiB | 7.89 GiB | 7.89 GiB | 7.90 GiB | 7.87 GiB | 7.88 GiB | 7.89 GiB | 7.89 GiB | 63.09 GiB |
| bytes_written_after_direct_io_disabled | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B |
| queue_size | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| work_units_consumed | 1873 | 1873 | 1874 | 1877 | 1870 | 1873 | 1873 | 1874 | 14987 |
| work_units_produced | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**reader:local**

| Key | Worker ID | | | | | | | | Total |
|---|---|---|---|---|---|---|---|---|---|
| | **0** | **1** | **2** | **3** | **4** | **5** | **6** | **7** | |
| bytes_consumed | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B |
| bytes_produced | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B |
| bytes_read | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B | 0 B |
| work_units_consumed | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| work_units_produced | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# ORGANIZING LOGS

/2016/06/14/frob_widgets_1.tar.bz2

# ORGANIZING LOGS

/2016/06/14/frob_widgets_1.tar.bz2

/cluster_nodes.txt

/node.conf

# ORGANIZING LOGS

/2016/06/14/frob_widgets_1.tar.bz2

/cluster_nodes.txt

/node.conf

/notes.md

SMARTLY CONTROLLING MEMORY

BUFFER POOLS ARE
FAST AND SIMPLE
AND INFLEXIBLE

# MALLOC IS
# SIMPLE AND FLEXIBLE
# AND DANGEROUS

By default, Linux follows an **optimistic** memory allocation strategy.  This means that when malloc() returns non-NULL there is **no guarantee that the memory really is available**.

In case it turns out that the system is out of memory, one or more processes **will be killed** by the OOM killer.

Apparently P has to be symmetric

$$\frac{1}{2}\vec{x}^T P \vec{x} + \vec{q}^T \vec{x} = c^2 - \sum_i q_i s_i + M^T \vec{x}$$

Very simple example: $\vec{x} = \begin{bmatrix} q_0 \\ s_0 \\ c \end{bmatrix}$

$$\frac{1}{2}\begin{bmatrix} q_0 & s_0 & c \end{bmatrix}\begin{bmatrix} P_{11} & P_{12} & P_{13} \\ P_{12} & P_{22} & P_{23} \\ P_{13} & P_{23} & P_{33} \end{bmatrix}\begin{bmatrix} q_0 \\ s_0 \\ c \end{bmatrix} = c^2 - \sum_i q_i s_i$$

$$\begin{bmatrix} P_{11}q_0 + P_{12}s_0 + P_{13}c, & P_{12}q_0 + P_{22}s_0 + \ldots \end{bmatrix}$$

$$\begin{bmatrix} \vec{x}^T \vec{P_1} & \vec{x}^T \vec{P_2} & \vec{x}^T \vec{P_3} \end{bmatrix}\begin{bmatrix} q_0 \\ s_0 \\ c \end{bmatrix} = c^2 - \sum_i q_i s_i$$

$$q_0 \vec{x}^T \vec{P_1} + \vec{x}^T \vec{P_2} s_0 + c \vec{x}^T \vec{P_3} = c^2 - q_0 s_0$$

$$q_0(q_0 P_{11} + s_0 P_{12} + c P_{13}) + s_0(q_0 P_{12} + s_0 P_{22} + c \ldots)$$
$$+ c(q_0 P_{13} + s_0 P_{23} + c P_{33}) = c^2 - q_0 s_0 - q_1 s_1$$

$$P_{11}q_0^2 + 2P_{12}q_0 s_0 + 2P_{23}c s_0 + 2P_{13}c q_0 + \ldots$$

$$= c^2 - q_0 s_0$$

(Note: lost ½ somewhere)
(mult by 2)

$P_{22} = 0$

$2P_{12} = -1$

$P_{33} = 1$

All else $\emptyset$

$$P = \begin{bmatrix} 0 & -\frac{1}{2} & 0 \\ -\frac{1}{2} & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} -s_0 & -q_0 & c \end{bmatrix}\begin{bmatrix} q_0 \\ s_0 \\ c \end{bmatrix} = \ldots$$

---

A →1 B →1 C
800  800  800
(3)  (2)  (1)

A →1 B →1 C
600  600  600
(3)  (1)  (2)

A →1 B 1:4 C
600  600  1200
(4)  (1)  (2)
     B → D
     600  1200
          (3)

$M_B = 4$

---

$\bar{C} \geq A$
$D \geq B$
$F \geq C + E + D$
$G \geq F$
$J \geq H + G$
$K \geq G + I$
$L \geq J$
$M \geq K$
$A = B$  $H = I$
$C = D$
$J = K$
$L = M$

$C + E + D \geq A + E + D$
$F \geq 2C + E$
$J \geq H + G$
$J \geq G + H$
$L \geq J$
$C \geq A$

---

$R_a M_a \geq \sum_{(a,y) \in E} P_{ay} R_y$

$R_a M_a \geq M_a \sum_{(x,a) \in E} P_{xa} R_x M_x$

$\sum_{(a,y) \in E} P_{ay} R_y \geq R_a M_a \geq \sum_{(x,a) \in E} P_{xa} R_a M_x$

Always holds even if you don't have incoming edges

$A \to B \to C$

$R_A \geq R_B$   $R_A \geq R_B \geq R_C$
$R_B \geq R_C$
$R_B \geq R_A$
$R_C \geq R_B$

$S \geq L + M$
$T \geq L + N$
Assume $I = M = N = \emptyset$
$S \geq L$
$T \geq L$
$J \geq G + H$
$L \geq J$

---

Upstream
$C \geq 0.5A + 0.5B$
$D \geq 0.5C$
$E \geq 0.5C$
Downstream
$A \geq 0.5C$
$B \geq 0.5C$
$C \geq D$
$C \geq E$

$A \geq C$   $C \geq A$
$B \geq C$   $C \geq B$
$C \geq 0.5D$   $D \geq 0.5C$
$C \geq 0.5E$   $E \geq 0.5C$

$P_{CD} = 1$
$P_{CD} \neq 0.5$

ARE WE SOLVING
THE RIGHT PROBLEM?

# WHAT IF MALLOC WAITED?

WAITING PROVIDES BACKPRESSURE

CALLERS CAN BE SCHEDULED

INTERFACE STAYS SIMPLE

DECISIONS CAN BE GLOBAL

# POOLS
# QUOTAS
# CONSTRAINTS

WRAPPING UP

# 938 GB PER MINUTE
# 15.6 GB PER SECOND
# 52 NODES
# 300 MBPS PER NODE

**6757** GB PER MINUTE

**112.6** GB PER SECOND

**178** NODES

**632** MBPS PER NODE

# LESSONS LEARNED:

# BOTTLENECKS SHAPE YOUR ARCHITECTURE

# STRUCTURE SOFTWARE FOR EXPERIMENTATION AND MEASUREMENT

SAVE YOUR LOGS
SAVE YOUR CONFIG
SAVE YOUR NOTES

SOMETIMES YOU NEED MORE CONTROL THAN THE OS WILL GIVE YOU

GOING FAST
IS HARD

# THANKS