# Cloud-hosted Data Transfer & Optimization:
# Stork for the Cloud

Tevfik Kosar

University at Buffalo (SUNY)

May 2, 2012

Condor Week, Madison, WI

# Big Data

## Science



- 1 PB is now considered "small" for many science applications today

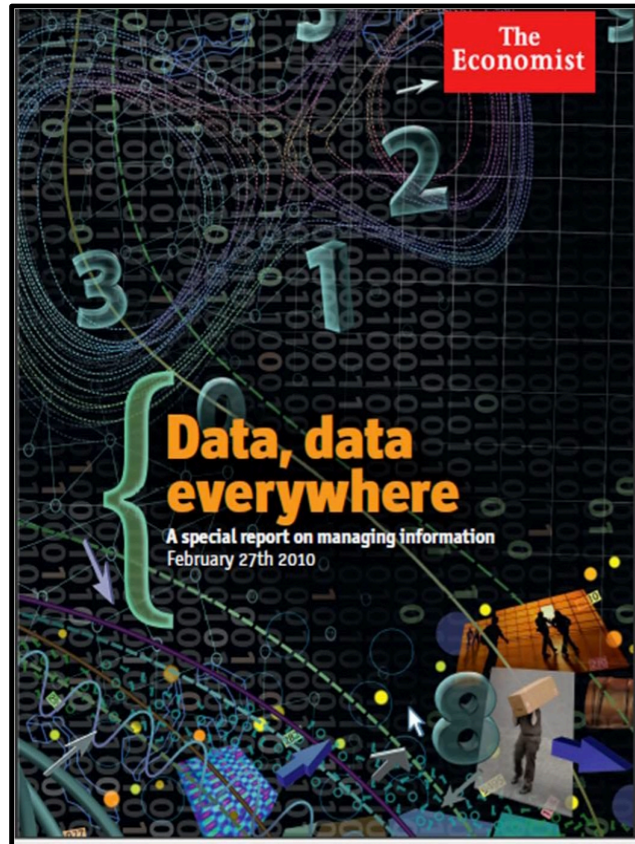- For most, their data is distributed across several sites

## Industry



A survey among 106 organizations operating two or more data centers:

- 77% run replication among three or more sites

- 50% has more than 1 PB in their primary data center

# Big Data

Industry



- 1 PB is now considered "small" for many science applications today

- For most, their data is distributed across several sites

A survey among 106 organizations operating two or more data centers:

- 77% run replication among three or more sites

- 50% has more than 1 PB in their primary data center

# Best Way to Move Big Data?

# Best Way to Move Big Data?

- Sending **1 PB** of data over 10 Gbps link would take **nine days** (assuming 100% efficiency) -- <span style="color:red">too optimistic!</span>
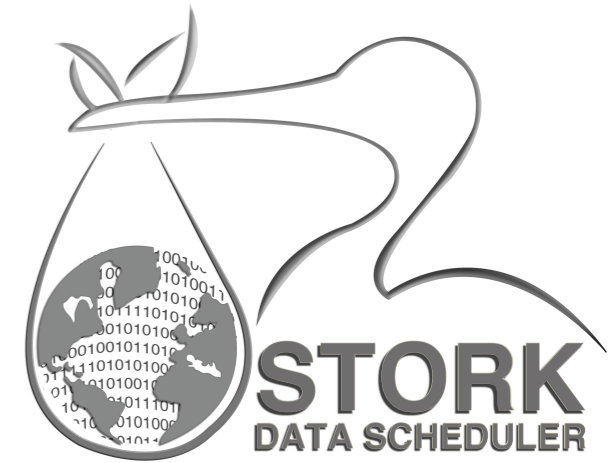
# Best Way to Move Big Data?

- Sending **1 PB** of data over 10 Gbps link would take **nine days** (assuming 100% efficiency) -- too optimistic!

- Sending **1 TB** Forensics dataset from Boston to Amazon S3 cost $100 and took **several weeks** [Garfinkel 2007]

# Best Way to Move Big Data?

- Sending **1 PB** of data over 10 Gbps link would take **nine days** (assuming 100% efficiency) -- too optimistic!

- Sending **1 TB** Forensics dataset from Boston to Amazon S3 cost $100 and took **several weeks** [Garfinkel 2007]

- Visualization scientists at LANL dumping data to tapes and sending them to Sandia Lab via **Fedex** [Feng 2003]

# Best Way to Move Big Data?

- Sending **1 PB** of data over 10 Gbps link would take **nine days** (assuming 100% efficiency) -- too optimistic!

- Sending **1 TB** Forensics dataset from Boston to Amazon S3 cost $100 and took **several weeks** [Garfinkel 2007]

- Visualization scientists at LANL dumping data to tapes and sending them to Sandia Lab via **Fedex** [Feng 2003]

- Collaborators have the option of moving their data into disks, and sending them as packages through **UPS or FedEx** [Cho et al 2011].

# Best Way to Move Big Data?

- Sending **1 PB** of data over 10 Gbps link would take **nine days** (assuming 100% efficiency) -- too optimistic!
- Sending **1 TB** Forensics dataset from Boston to Amazon S3 cost $100 and took **several weeks** [Garfinkel 2007]
- Visualization scientists at LANL dumping data to tapes and sending them to Sandia Lab via **Fedex** [Feng 2003]
- Collaborators have the option of moving their data into disks, and sending them as packages through **UPS or FedEx** [Cho et al 2011].
- Will **100 Gbps** networks change anything?

# Stork Data Scheduler

- Implements state-of-the art models and algorithms for data scheduling & optimization

- Started as part of the Condor Project (was my PhD work)

- Currently developed at University at Buffalo and funded by NSF (CAREER, STCI, CiC)

- Based on the Condor code, uses Condor libraries (DaemonCore, ClassAds)

- Compatible with Condor products (i.e. DAGMan)

- .....

# Stork Data Scheduler

- .....

- Built & tested on Condor NMI (Metronome)

- Supports more than 20 platforms

- Futures include:
    - support for multiple transfer protocols
    - dynamic protocol tuning & optimization
    - end-to-end throughput prediction services
    - data aggregation & connection caching
    - early error detection and classification & recovery

# End-to-end Problem

Data flow
Control flow

CPU    NIC    **Tnetwork**    NIC    CPU

**TDnetwork->mem**

Memory    **TSmem->network**    Memory

**TSdisk->mem**    **TDmem->disk**

DISK    DISK

**Tnetwork** -> Network Throughput
**TSmem->network** -> Memory-to-network
Throughput on source
**TSdisk->mem** -> Disk-to-memory Throughput on
source
**TDnetwork->mem** -> Network-to-memory
Throughput on Destination
**TDmem->disk** -> Memory-to-disk Throughput on
destination

# End-to-end Problem

Data flow

Control flow

CPU    NIC

**Tnetwork**

NIC    CPU

**TSmem->network**

**TDnetwork->mem**

Memory

Memory

**TSdisk->mem**

**TDmem->disk**

DISK

DISK

**Tnetwork** -> Network Throughput
**TSmem->network** -> Memory-to-network
Throughput on source
**TSdisk->mem** -> Disk-to-memory Throughput on
source
**TDnetwork->mem** -> Network-to-memory
Throughput on Destination
**TDmem->disk** -> Memory-to-disk Throughput on
destination

protocol
tuning

# End-to-end Problem

Data flow ──────
Control flow ···········



**Tnetwork**

**TSmem->network**

**TDnetwork->mem**

**TSdisk->mem**

**TDmem->disk**

**Tnetwork** -> Network Throughput
**TSmem->network** -> Memory-to-network Throughput on source
**TSdisk->mem** -> Disk-to-memory Throughput on source
**TDnetwork->mem** -> Network-to-memory Throughput on Destination
**TDmem->disk** -> Memory-to-disk Throughput on destination

protocol
tuning



disk I/O
optimization

# End-to-end Problem



Data flow
Control flow

CPU   NIC   **Tnetwork**   NIC   CPU

**TSmem->network**   **TDnetwork->mem**

Memory   Memory

**TSdisk->mem**   **TDmem->disk**

DISK   DISK

**Tnetwork** -> Network Throughput
**TSmem->network** -> Memory-to-network
Throughput on source
**TSdisk->mem** -> Disk-to-memory Throughput on
source
**TDnetwork->mem** -> Network-to-memory
Throughput on Destination
**TDmem->disk** -> Memory-to-disk Throughput on
destination

protocol
tuning

disk I/O
optimization

CPU
optimization

# End-to-end Problem



Data flow
Control flow

**Tnetwork** -> Network Throughput
**TSmem->network** -> Memory-to-network Throughput on source
**TSdisk->mem** -> Disk-to-memory Throughput on source
**TDnetwork->mem** -> Network-to-memory Throughput on Destination
**TDmem->disk** -> Memory-to-disk Throughput on destination

protocol tuning

disk I/O optimization

CPU optimization

Parameters to be optimized:
- # of streams
- # of disk stripes
- # of CPUs/nodes

# End-to-end Optimization



- CPU nodes are considered as nodes of a maximum flow problem

- Memory-to-memory transfers are simulated with dummy source and sink nodes

- The capacities of disk and network is found by applying parallel stream model by taking into consideration of resource capacities (NIC & CPU)

# Challenging Problem

Optimize:

- concurrency
- parallelism
- pipelining
- conn. caching
- buffer size
- block size
- disk striping
- threading
- ....

# Challenging Problem

(1)

Optimize:

- concurrency
- parallelism
- pipelining
- conn. caching
- buffer size
- block size
- disk striping
- threading
- ....



512 x 8 MB files

# Challenging Problem

Optimize:

- concurrency
- parallelism
- pipelining
- conn. caching
- buffer size
- block size
- disk striping
- threading
- ....

(1)



512 x 8 MB files

# Challenging Problem

Optimize:

- concurrency
- parallelism
- pipelining
- conn. caching
- buffer size
- block size
- disk striping
- threading
- ....

(1)



512 x 8 MB files

# Challenging Problem

(2)

Optimize:

- concurrency
- parallelism
- pipelining
- conn. caching
- buffer size
- block size
- disk striping
- threading
- ....



512 x 1 MB files

# Challenging Problem

(2)

Optimize:

- concurrency
- parallelism
- pipelining
- conn. caching
- buffer size
- block size
- disk striping
- threading
- ....



512 x 1 MB files

# Challenging Problem

Optimize:

- concurrency
- parallelism
- pipelining
- conn. caching
- buffer size
- block size
- disk striping
- threading
- ....



512 x 32 MB files

# Challenging Problem

(3)

Optimize:

- concurrency
- parallelism
- pipelining
- conn. caching
- buffer size
- block size
- disk striping
- threading
- ....



512 x 32 MB files

# Kosar et al Models



Exponential Packet Loss

Break Function Modeling

Modeling Based on Newton's Iteration

$$p'_n = a'n^{c'} + b'$$

Modeling Based on Full Second Order

$$p'_n = p_n \frac{RTT_n^2}{c^2 MSS^2} = a'n^2 + b'n + c'$$

# Kosar et al Models

- Details in 2 TPDS 2011 papers

- Implemented in the latest version of Stork (v.2.0.1)

- Provides throughput optimization as well as estimation

Modeling Based on Newton's Iteration

$$p'_n = a' n^{c'} + b'$$

Modeling Based on Full Second Order

$$p'_n = p_n \frac{RTT_n^2}{c^2 MSS^2} = a' n^2 + b' n + c'$$

# Stork for the Cloud

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork Android Client

# Stork for the Cloud



- Prototype implementation complete, testing stage
- Will be deployed as hosted service
- Allow deployment on private clouds as well
- Available on Amazon EC2 and Windows Azure
- More optimizations coming

# 100 Gbit Performance

# Summary

- Scientific and commercial applications are getting more and more data intensive

- Data sharing and bulk data transfers are still a  major bottleneck in front of multi-institutional and inter-disciplinary collaborative science

- Stork for the Cloud provides end-to-end throughput optimization in hosted environment accessible through ultra-thin clients

CYBERINFRASTRUCTURE VISION
FOR 21ST CENTURY DISCOVERY

National Science Foundation
Cyberinfrastructure Council
March 2007

The
FOURTH
PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

This work has been sponsored by:
# NSF, DOE, ONR, NOAA

For more information:
**Stork web page**: http://www.storkproject.org

**Questions?**

This work has been sponsored by:
**NSF, DOE, ONR, NOAA**

For more information:
**Stork web page**: http://www.storkproject.org