# Data Storage Roundtable

**Steven Ludtke**, Session Chair
Professor, Department of Biochemistry and Molecular Biology
Baylor College of Medicine, Houston

**Grant Jensen**
Professor, Department of Biology
California Institute of Technology, Pasadena

**David Mastronarde**
Associate Professor Department of MCD Biology
University of Colorado, Boulder

**Roberto Marabini**
Escuela Superior Politécnica Superior
Universidad Autónoma de Madrid, SPAIN

**Ardan Patwardhan**
Senior Scientific Officer, Protein Data Bank in Europe
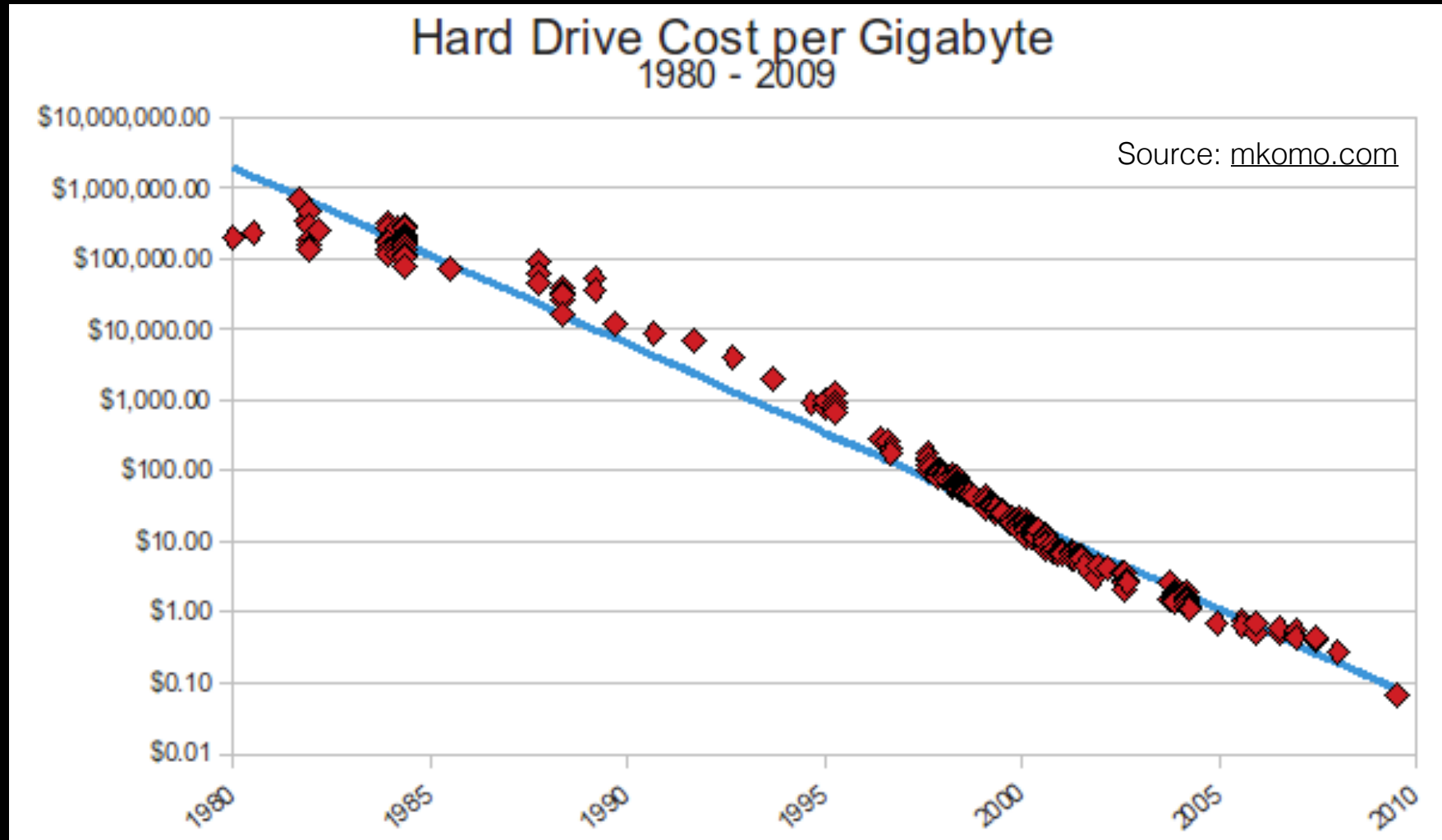European Bioinformatics Institute, Cambridge, UK

# Storage Issues

- Quantity of Data (10 TB - 10 PB)

- Data Bandwidth

- Reliability/Redundancy

- Cost

- Tomography vs High Resolution Movies

- Central archives/databases

# Quantity of Data

- 8x x 8k super-resolution counting movie, 30 frames

- 2 gigapixels/per movie

- typically only a few counts per pixel per frame

- 4 bits -> 1 GB movie (plus 256 MB periodically)

- 32 bits -> 8 GB movie

- Compression (slow, but saves even more space)

- Krios+K2 assume 0.5 - 1 TB/day

Ludtke, 4/2016

# Cost Over Time



Hard Drive Cost per Gigabyte
1980 - 2009

Source: mkomo.com

On average storage cost falls 2x every 14 months !
Most enterprise drives have 5 year warranty

Ludtke, 4/2016

# How Much Speed do You Need ?

- Xeon E5-2697v2
  - ~500 GFlops
  - ~200 GOps
  - 100 GB/sec memory bandwidth

- @100 MB/sec:
  - 5000 Flops/byte
  - 2000 Ops/byte

If a job processes 10 GB of data and takes 1 hour to run, should you worry about I/O speed ?

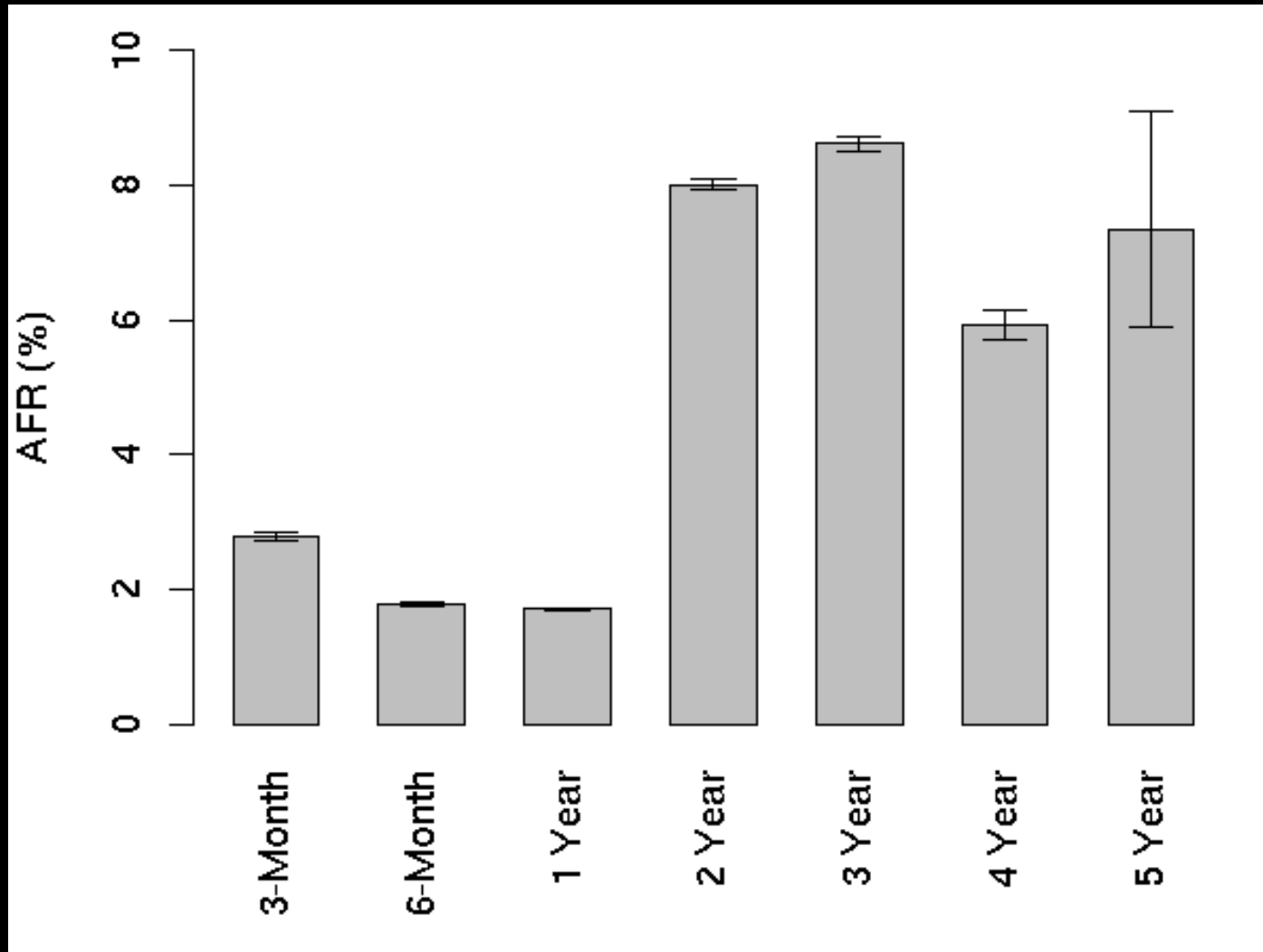How about a job where processing 10 GB of data takes only 10 seconds ?

# Interface Data Bandwidth

| | Speed (GB/sec) | Time to Transfer 4 TB |
|---|---|---|
| USB2 | 0.04 | 28 hours |
| Gigabit Network | 0.1 (0.125) | 11 hours |
| USB3 | ~0.3 | 3-4 hours |
| SATA | 0.3, 0.6, 1.2 | 1-3 hours |
| 10Gb Network | 1.0 (1.25) | 1 hour |
| Thunderbolt 2 | 2.0 | 30 min |
| Infiniband | 1.0-4.0 | 15 min - 1 hour |
| PCIe 3.0 | ~1.0/channel ~16.0 max | ~4 min |

# Drive Data Bandwidth

| | Speed (GB/sec) | Max Size |
|---|---|---|
| 2.5" Spinning Platter | 0.06 | 2 TB |
| 3.5" Spinning Platter | 0.1-0.2 | 8 TB |
| 2.5" SSD | 0.3-0.6 | 2 TB (16 TB) |
| RAID (striping) | 1.0-1.5 Typ 3.0+ Possible | ~80 TB/Array |
| PCIe/m.2 SSD ($$$) | ~2.0 Typ 4.0-6.0 Possible | ~1 TB |

# Annual Failure Rates

Ludtke, 4/2016

# Backup Concepts

- Offsite !

  - In case of physical disaster (hurricane, flood,…)

  - Have to arrange for space in another facility

  - Bandwidth available ?

- Duplicate hardware with 2nd copy

  - "Batch" problems with hardware

  - Hackers (intentional destruction of data)

  - Double the cost

- "offline" storage

  - Drives on a shelf - Human effort & "exercising"

  - Tape libraries - Human effort or robot ?   2.5 TB tape ~$80

# 56 TB - an Example

Workstation with 8-bay chassis + PCIe RAID controller

Cost w 5/3 year warranty ~$4000 —> $1.20/TB-month (+comp)
~1.3 GB/sec, and is also a computer!

Workstation cost +~$10,000
28 Cores, 2.6 Ghz, 128 GB RAM, GPU

(Note that this machine can be cheaper.
This configuration permits up to 4 GPUs.
Beware companies that sell or lease you 'threads' or
'virtual cores' as 'cores'. NOT the same!)

Advantages: Fast, movie processing!
Disadvantages: Expandability

# 80 TB - an Example

12 Bay Synology - $1300
12x 8TB He8 Drives - $6000

RAID6 -> ~80TB
Cost w 5 year warranty:
    ~$7300 —> $1.52/TB-month

~0.1 GB/sec (network limited)



Advantages: Reliable, Easy, Quiet, Cheap
Disadvantages: Slow

# 540 TB - an Example

1x4U computer with 36x 8TB drives ($24,000)+
1x4U 44x 8TB drives JBOD Chassis ($26,000)
Configured as 6x RAID6 volumes —> 540 TB usable
~1.5 GB/sec I/O to the attached computer

Cost w 3/5 year warranty ~$50k —> $1.54/TB-month
x5 —> 2.7 PB/rack (usable)



Advantages: Inexpensive, Fast, Includes Computing
Disadvantages: Management, Housing/Noise

Ludtke, 4/2016

# Cloud Storage ?

Amazon (S3):
- Standard Storage: $29.50/TB-month
- Infrequent Access: $12.50/TB-month
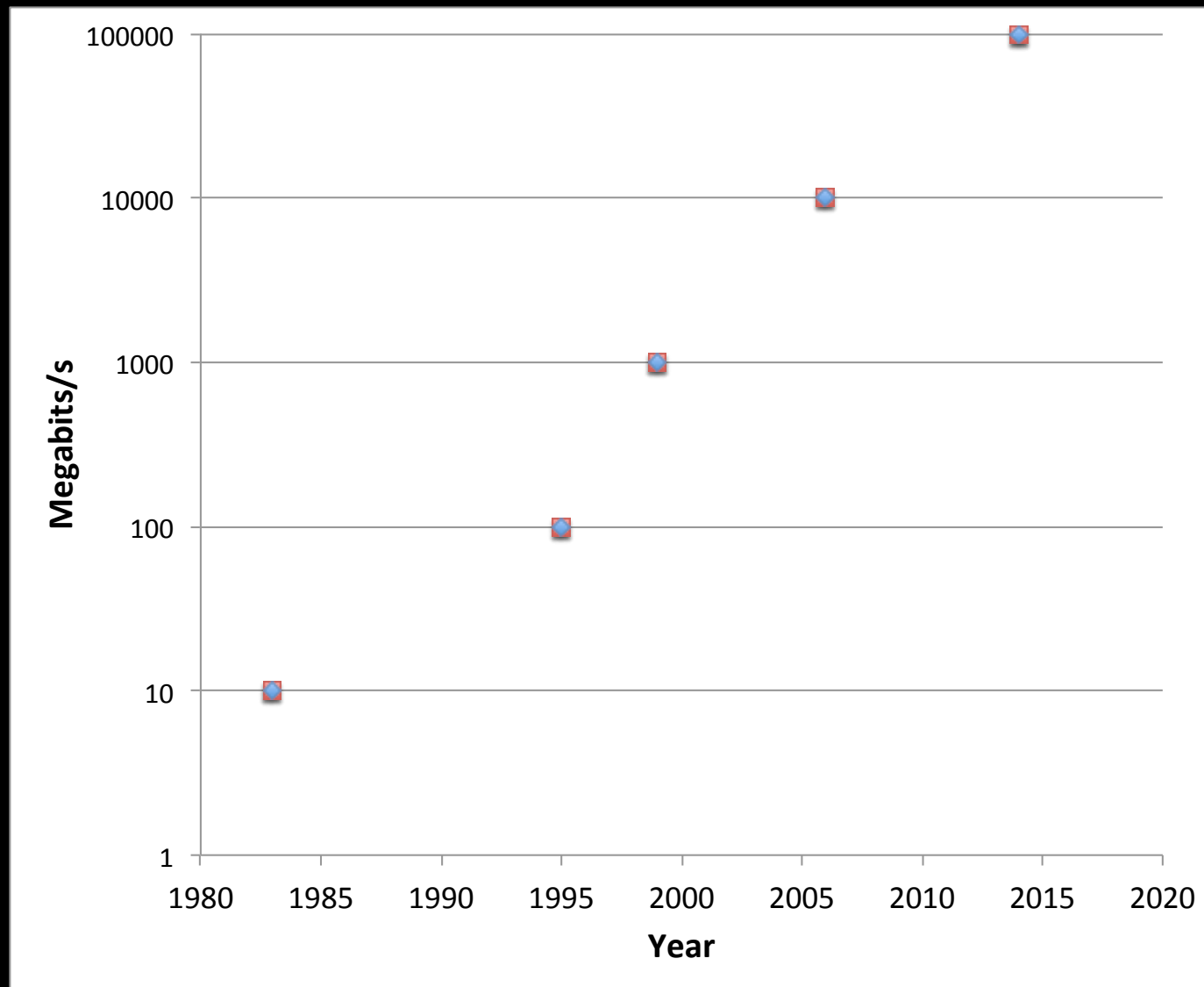- Glacier Storage (backup): $7/TB-month

+

Download cost:
- $50/TB

Advantages: Safe & Reliable, Access to EC2
Disadvantages: Slow Access, Expensive, Legal Issues

# Network Bandwidth



On average network bandwidth doubles every 27 months
Capacity doubles every 14 months!

Ludtke, 4/2016

# Hidden Costs ?

With second unit for backup —> $100k (total)

Administration costs
- Sysadmin $60-80k/year
- Amazon storage also needs to be locally managed!

Housing Costs (?)
- Depends on circumstances
- 1/5 Rack @coloc ~$20k for 5 years

Fractional Usage
- If you buy all at once, but gradually fill over lifetime, effective cost goes up