

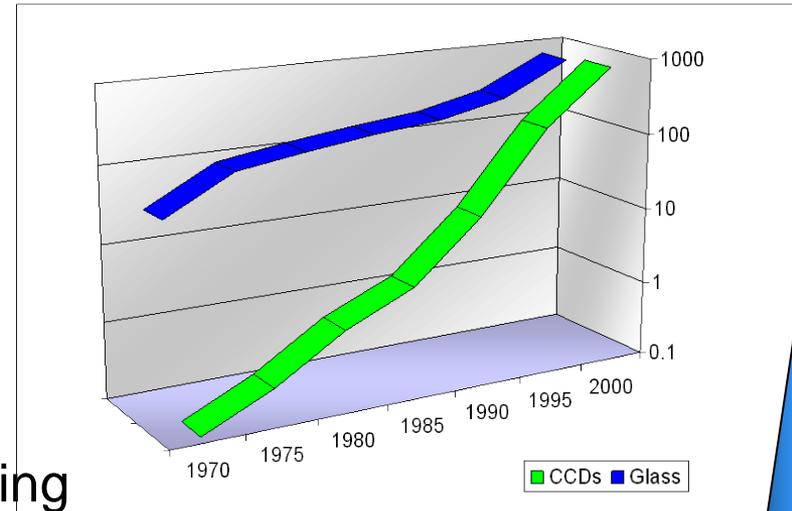


# Amdahl's Laws and Extreme Data-Intensive Computing

Alex Szalay  
The Johns Hopkins University

# Living in an Exponential World

- Scientific data doubles every year
  - *caused by successive generations of inexpensive sensors + exponentially faster computing*
- Changes the nature of scientific computing
- Cuts across disciplines (eScience)
- It becomes increasingly harder to extract knowledge
- 20% of the world's servers go into huge data centers by the “Big 5”
  - *Google, Microsoft, Yahoo, Amazon, eBay*



# Collecting Data

- Very extended distribution of data sets:  
***data on all scales!***
- Most datasets are small, and manually maintained (Excel spreadsheets)
- Total amount of data dominated by the other end (large multi-TB archive facilities)
- Most bytes today are collected via electronic sensors



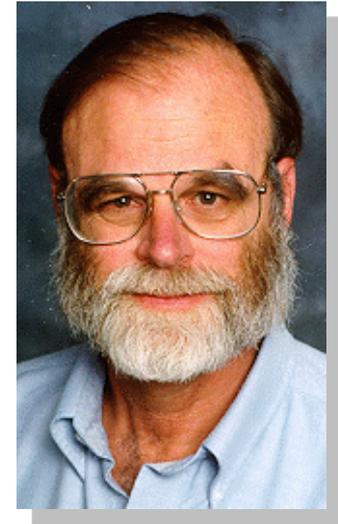
# Scientific Data Analysis

- Data is everywhere, never will be at a single location
- Architectures increasingly CPU-heavy, IO-poor
- Data-intensive scalable architectures needed
- Need randomized, incremental algorithms
  - *Best result in 1 min, 1 hour, 1 day, 1 week*
- Most scientific data analysis done on small to midsize BeoWulf clusters, from faculty startup
- Universities hitting the “power wall”
- **Not scalable, not maintainable...**

# Gray's Laws of Data Engineering

## Jim Gray:

- Scientific computing is revolving around **data**
- Need **scale-out** solution for analysis
- Take the **analysis to the data!**
- Start with “**20 queries**”
- Go from “**working to working**”



DISC: Data Intensive Scientific Computing

# Building Scientific Databases

- 10 years ago we set out to explore how to cope with the data explosion (with Jim Gray)
- Started in astronomy, with the Sloan Digital Sky Survey
- Expanded into other areas, while exploring what can be transferred
- During this time data sets grew from 100GB to 100TB
- Interactions with every step of the scientific process
  - *Data collection, data cleaning, data archiving, data organization, data publishing, mirroring, data distribution, data curation...*

# Reference Applications

## Some key projects at JHU

- **SDSS:** *100TB total, 35TB in DB, in use for 8 years*
- **NVO :** *~5TB, few billion rows, in use for 4 years*
- **PanStarrs:** *80TB by 2011, 300+ TB by 2012*
- **Immersive Turbulence:** *30TB now, 100TB by Dec 2010*
- **Sensor Networks:** *200M measurements now, forming complex relationships*

## Key Questions:

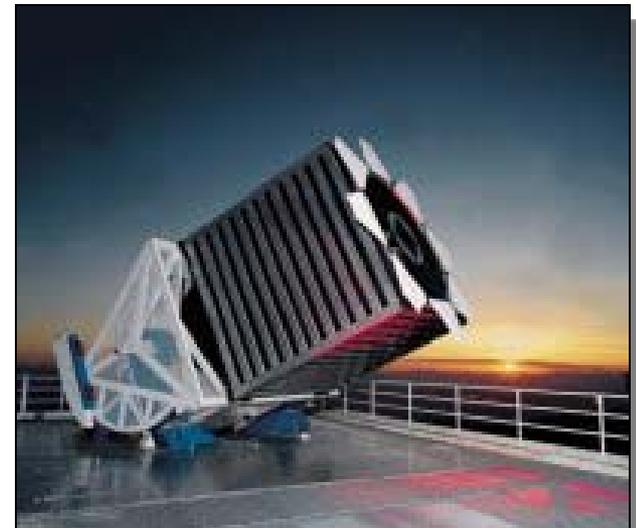
- What are the reasonable tradeoffs for DISC?
- How do we build a ‘scalable’ architecture?
- How do we interact with petabytes of data?

# Sloan Digital Sky Survey



- “The Cosmic Genome Project”
- Two surveys in one
  - Photometric survey in 5 bands
  - Spectroscopic redshift survey
- Data is public
  - 40 TB of raw data
  - 5 TB processed catalogs
  - 2.5 Terapixels of images
- Started in 1992, finishing in 2008
- Database and spectrograph built at JHU (SkyServer)

*The University of Chicago  
Princeton University  
The Johns Hopkins University  
The University of Washington  
New Mexico State University  
Fermi National Accelerator Laboratory  
US Naval Observatory  
The Japanese Participation Group  
The Institute for Advanced Study  
Max Planck Inst, Heidelberg  
Sloan Foundation, NSF, DOE, NASA*

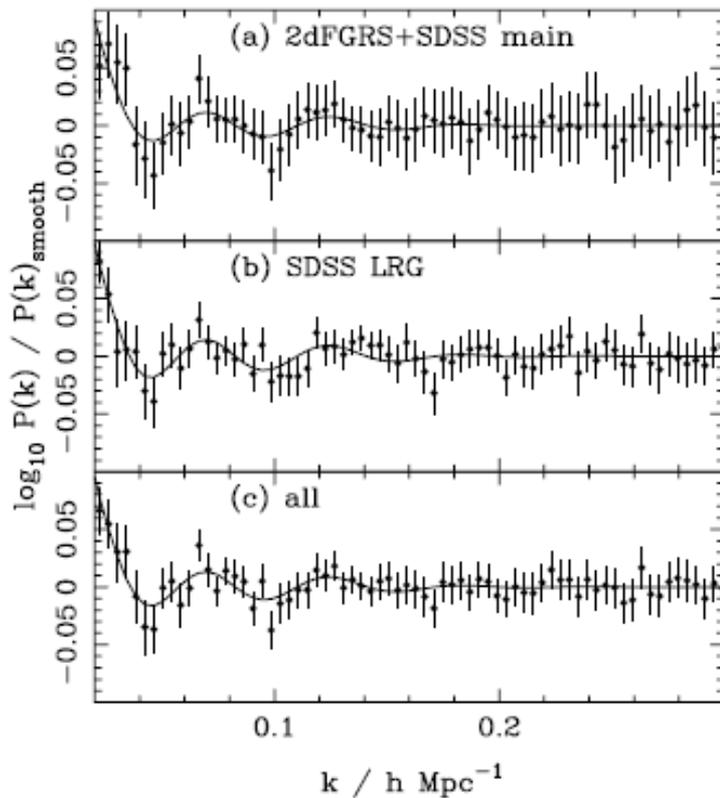


# Primordial Sound Waves in SDSS

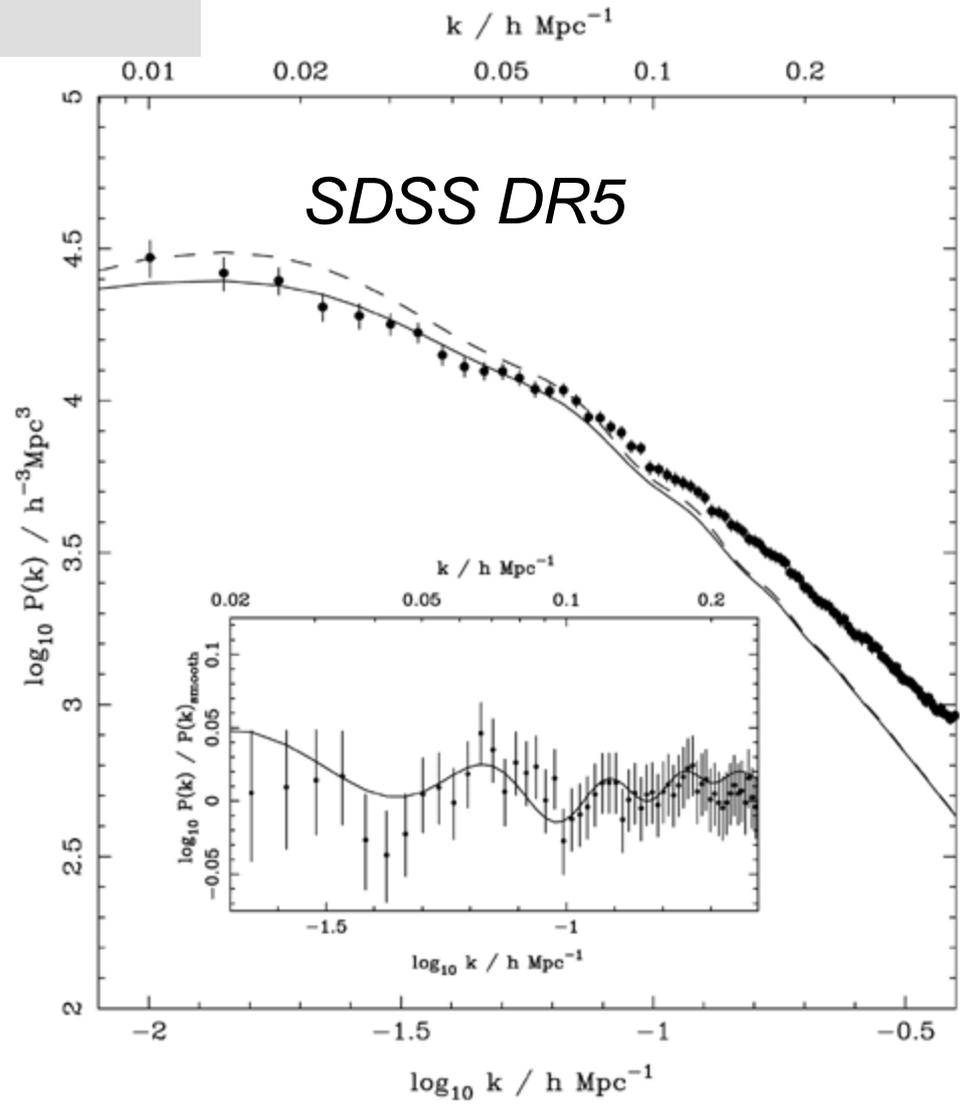
800K galaxies

Power Spectrum

(Percival et al 2006, 2007)



*SDSS DR6+2dF*



# Public Use of the SkyServer

- **Prototype in 21<sup>st</sup> Century**

- 400 million web hits in 6 years
- 930,000 distinct users
- 100 million images vs 10,000 astronomers
- Delivered 50,000 hours of lectures to high schools
- Delivered 100B rows of data
- Everything is a power law



- **GalaxyZoo**

- 40 million visual galaxy classifications by the public
- Enormous publicity (CNN, Times, Washington Post, BBC)
- 100,000 people participating, blogs, poems, ....
- Now truly amazing original discovery by a schoolteacher

# Pan-STARRS



- **Detect 'killer asteroids'**
  - *PS1: starting in May 1, 2010*
  - *Hawaii + JHU + Harvard/CfA + Edinburgh/Durham/Belfast + Max Planck Society*
- **Data Volume**
  - *>1 Petabytes/year raw data*
  - *Camera with 1.4Gigapixels*
  - *Over 3B celestial objects plus 250B detections in database*
  - *80TB SQLServer database built at JHU, 3 copies for redundancy*

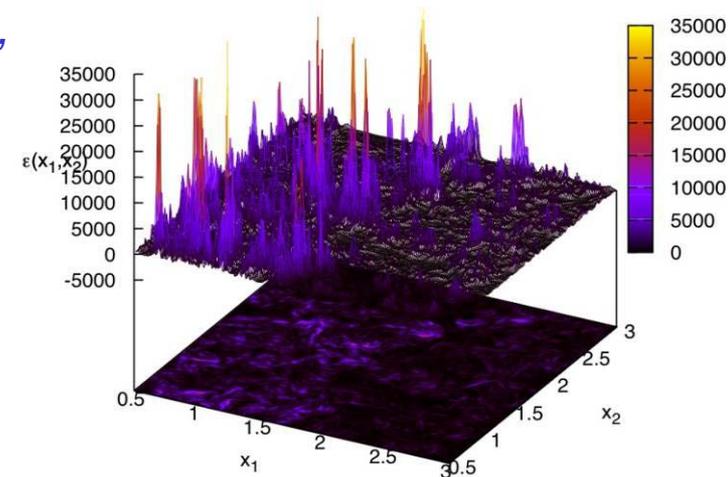
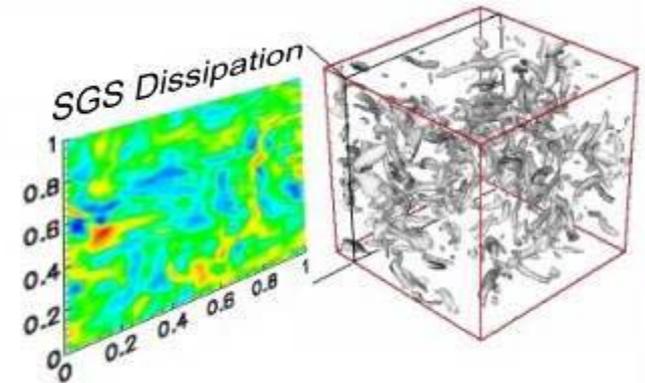


# Immersive Turbulence

- **Understand the nature of turbulence**

- *Consecutive snapshots of a  $1,024^3$  simulation of turbulence: now 30 Terabytes*
- *Treat it as an experiment, observe the database!*
- *Throw test particles (sensors) in from your laptop, immerse into the simulation, like in the movie Twister*

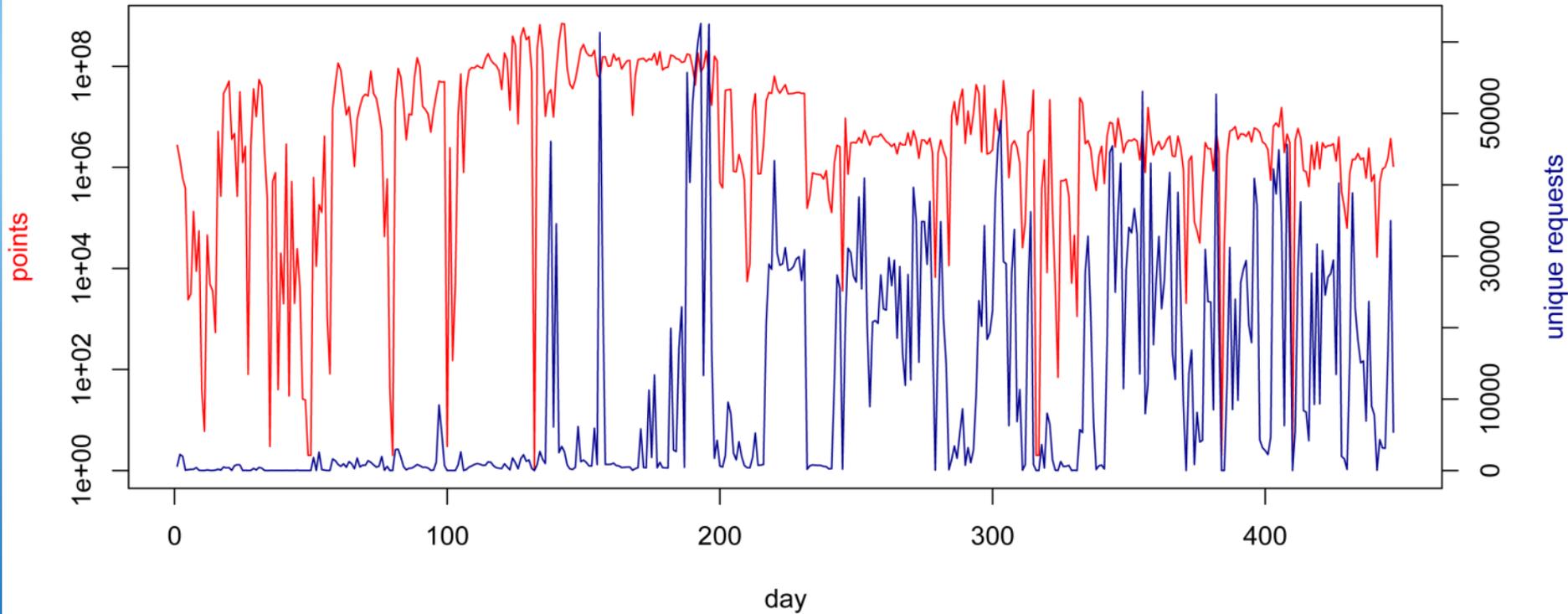
- **New paradigm** for analyzing HPC simulations!



with C. Meneveau, S. Chen (Mech. E), G. Eyink (Applied Math), R. Burns (CS)

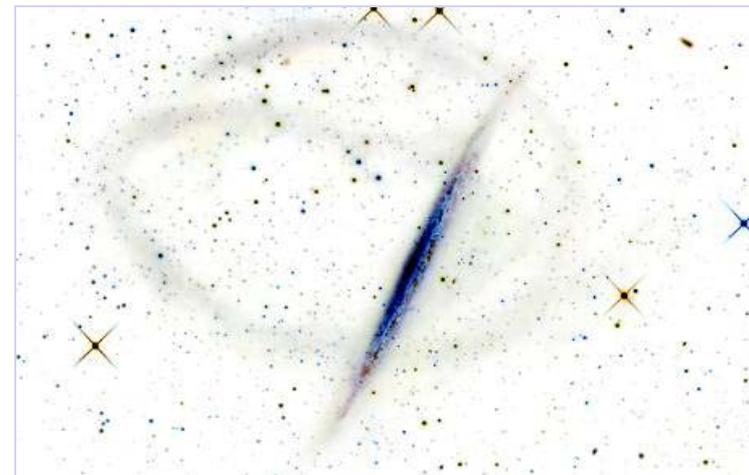
# Daily Usage

Turbulence Database Usage by Day



# The Milky Way Laboratory

- Pending NSF Proposal to use cosmology simulations as an immersive laboratory for general users
- Use Via Lactea-II (20TB) as prototype, then Silver River (500TB+) as production (15M CPU hours)
- Output 10K+ hi-rez snapshots (200x of previous)
- Users insert test particles (dwarf galaxies) into system and follow trajectories in precomputed simulation
- Realistic “streams” from tidal disruption
- Users interact remotely with 0.5PB in ‘real time’



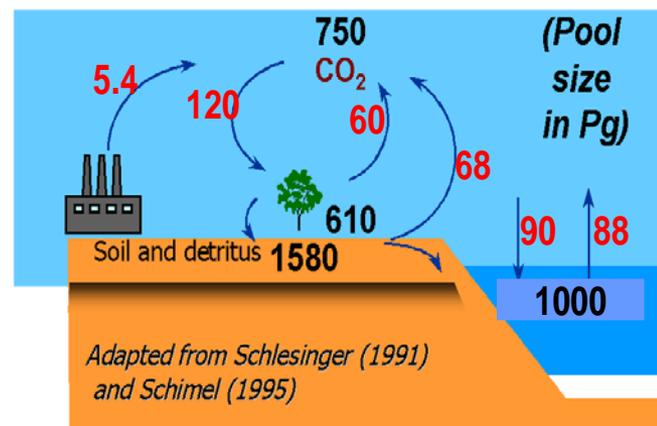
# Life Under Your Feet

- **Role of the soil in Global Change**

- *Soil CO<sub>2</sub> emission thought to be >15 times of anthropogenic*
- *Using sensors we can measure it directly, in situ, over a large area*

- **Wireless sensor network**

- *Use 100+ wireless computers (motes), with 10 sensors each, monitoring*
  - *Air +soil temperature, soil moisture, ...*
  - *Few sensors measure CO<sub>2</sub> concentration*
- *Long-term continuous data, 180K sensor days, 30M samples*
- *Complex database of sensor data, built from the SkyServer*
- *End-to-end data system, with inventory and calibration databases*

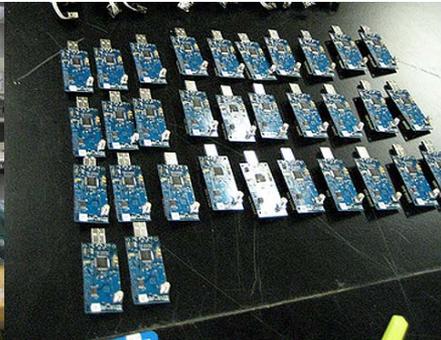


with K.Szlavec (Earth and Planetary), A. Terzis (CS)

<http://lifeunderyourfeet.org/>

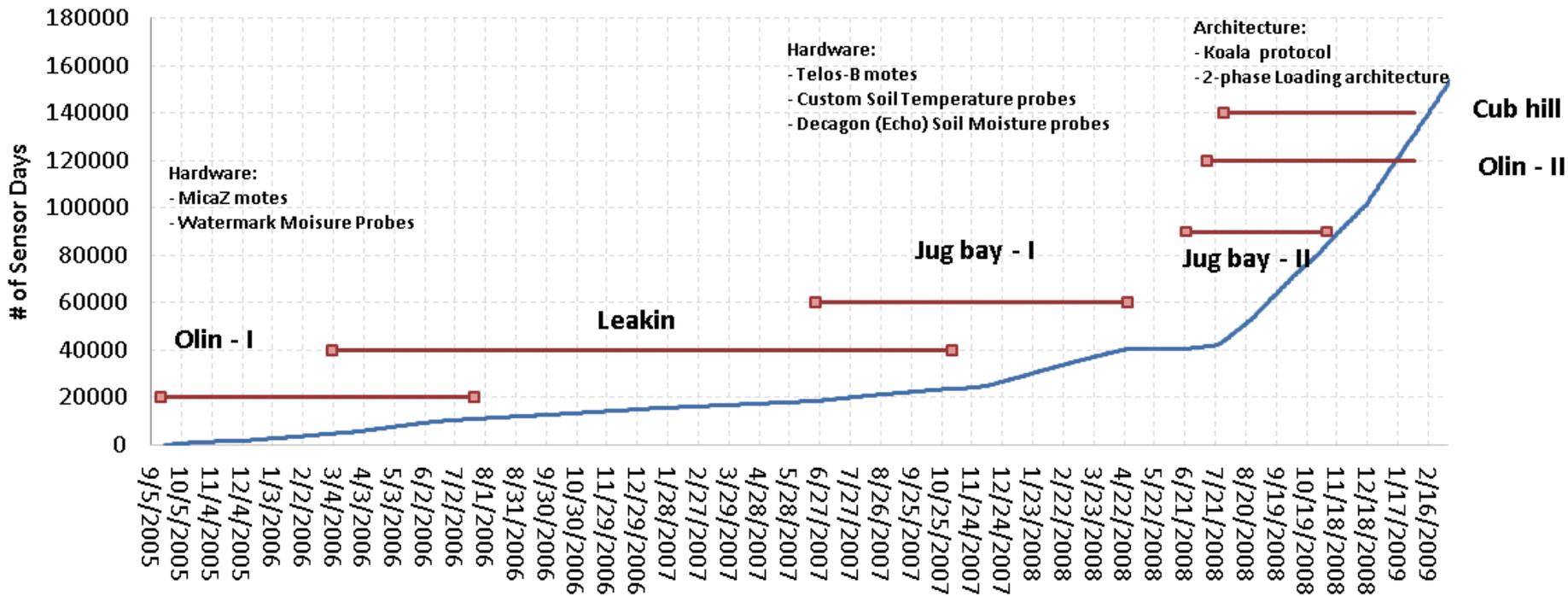
# Current Status

- Designed and built 2nd generation mote platform
  - *Telos SkyMote + own DAQ board*
- Hierarchical network architecture (Koala)
- Improved mote software
  - *Support for large-scale deployments*
  - *Over-the-air reprogramming*
  - *Daily log file written*
  - *Increased power efficiency (2 years on a single battery)*



# Cumulative Sensor Days

LUYF Sensor days



# Commonalities

- Huge amounts of data, aggregates needed
  - *But also need to keep raw data*
  - *Need for parallelism*
- Use patterns enormously benefit from indexing
  - *Rapidly extract small subsets of large data sets*
  - *Geospatial everywhere*
  - *Compute aggregates*
  - *Fast sequential read performance is critical!!!*
  - *But, in the end everything goes.... search for the unknown!!*
- Data will never be in one place
  - *Newest (and biggest) data are live, changing daily*
- Fits DB quite well, but no need for transactions
- Design pattern: class libraries wrapped in SQL UDF
  - *Take analysis to the data!!*

# Continuing Growth

## How long does the data growth continue?

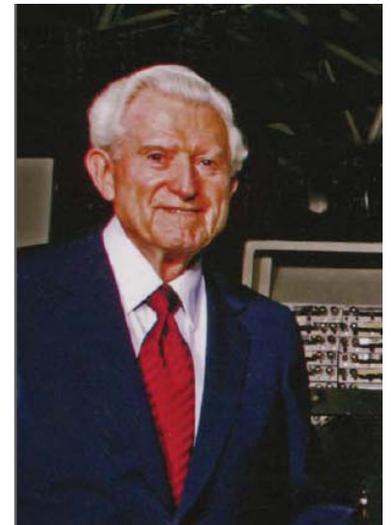
- High end always linear
- Exponential comes from technology + economics
  - ↔ rapidly changing generations
    - *like CCD's replacing plates, and become ever cheaper*
- How many new generations of instruments do we have left?
- Are there new growth areas emerging?
- **Software is becoming a new kind instrument**
  - *Value added federated data sets*
  - *Simulations*
  - *Hierarchical data replication*

# Amdahl's Laws

Gene Amdahl (1965): Laws for a balanced system

- i. Parallelism: max speedup is  $S/(S+P)$
- ii. **One bit of IO/sec per instruction/sec (BW)**
- iii. **One byte of memory per one instruction/sec (MEM)**

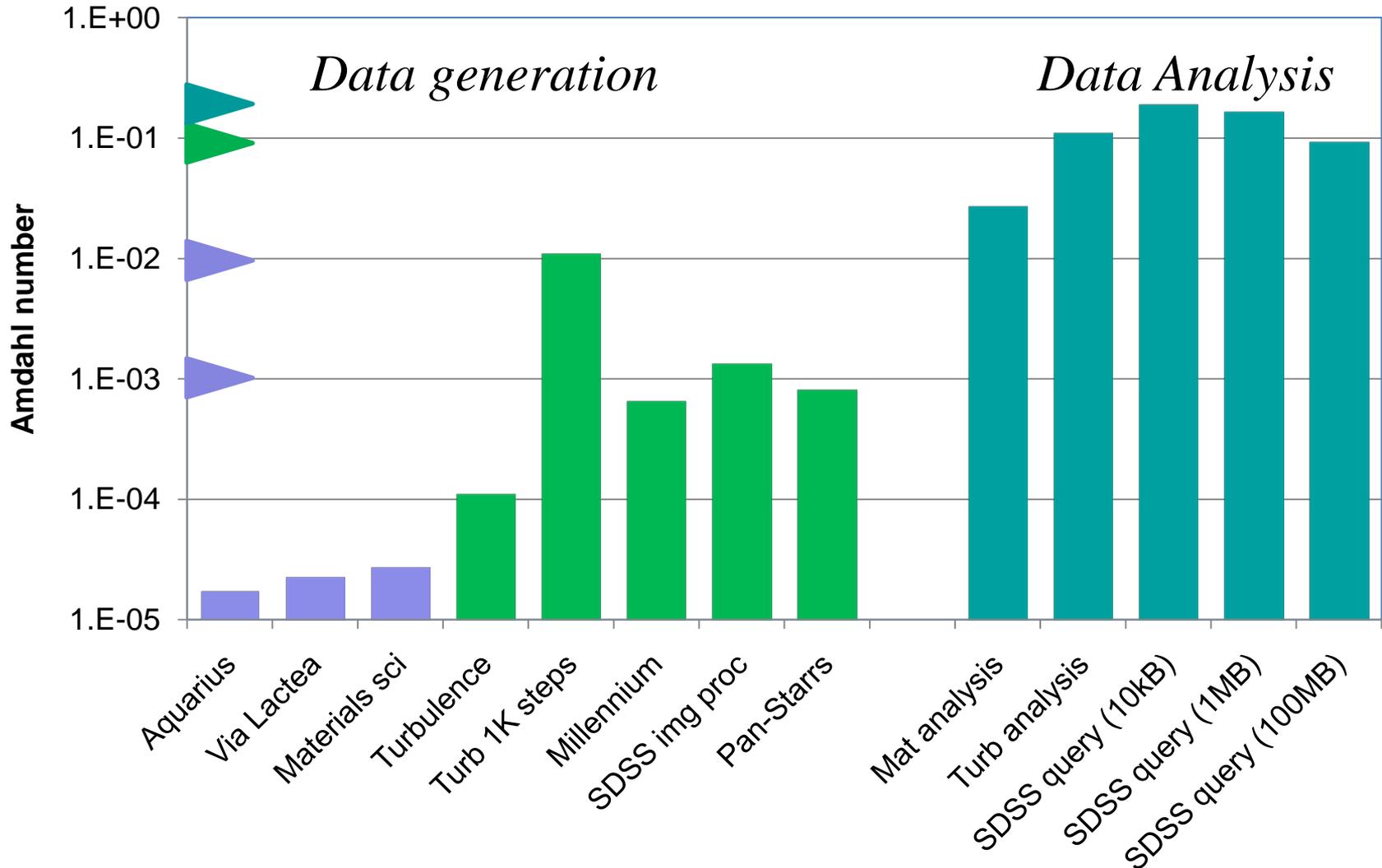
Modern multi-core systems move farther  
away from Amdahl's Laws  
(Bell, Gray and Szalay 2006)



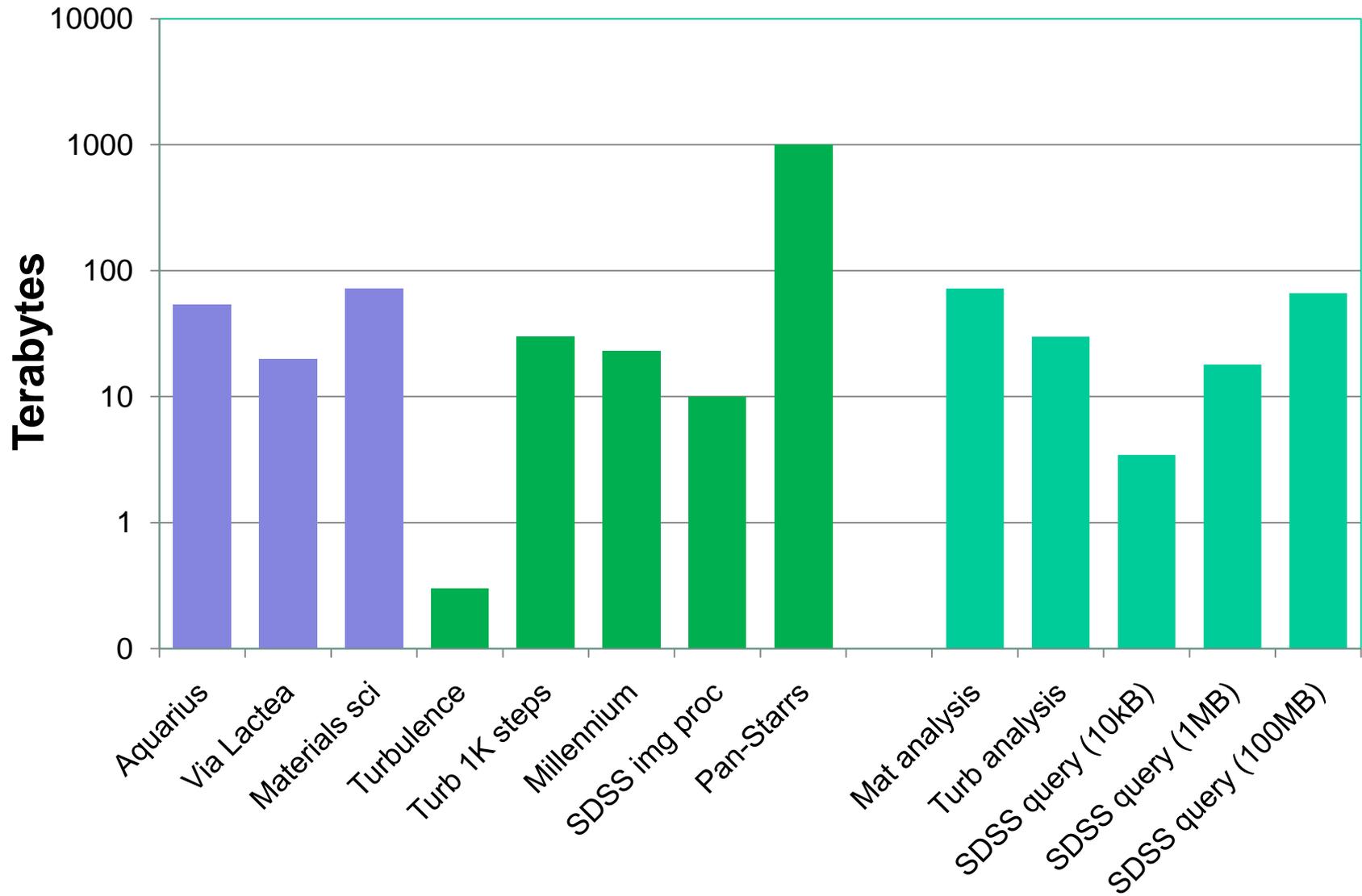
# Typical Amdahl Numbers

<i>System</i>	<i>CPU count</i>	<i>GIPS [GHz]</i>	<i>RAM [GB]</i>	<i>diskIO [MB/s]</i>	<i>Amdahl</i>	
					<i>RAM</i>	<i>IO</i>
<i>BeoWulf</i>	100	300	200	3000	0.67	0.08
<i>Desktop</i>	2	6	4	150	0.67	0.2
<i>Cloud VM</i>	1	3	4	30	1.33	0.08
<i>SC1</i>	212992	150000	18600	16900	0.12	0.001
<i>SC2</i>	2090	5000	8260	4700	1.65	0.008
<i>GrayWulf</i>	416	1107	1152	70000	1.04	0.506

# Amdahl Numbers for Data Sets



# The Data Sizes Involved

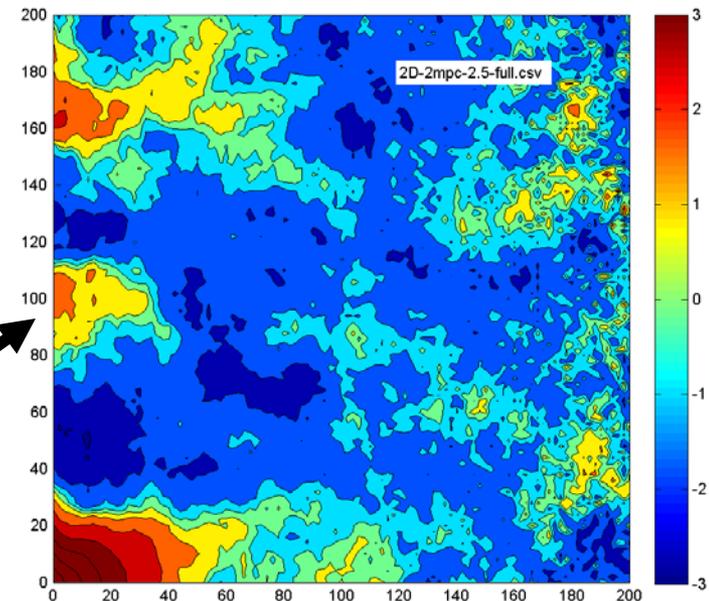


# DISC Needs Today

- Disk space, disk space, disk space!!!!
- Current problems not on Google scale yet:
  - *10-30TB easy, 100TB doable, 300TB really hard*
  - *For detailed analysis we need to park data for several mo*
- Sequential IO bandwidth
  - *If not sequential for large data set, we cannot do it*
- How do can move 100TB within a University?
  - *1Gbps                      10 days*
  - *10 Gbps                      1 day (but need to share backbone)*
  - *100 lbs box                      few hours*
- From outside?
  - *Dedicated 10Gbps or FedEx*

# The Impact of GPUs

- We need to reconsider the  $N \log N$  only approach
- Once we can run 100K threads, maybe running SIMD  $N^2$  on smaller partitions is also acceptable
- Potential impact on genomics huge
  - *Sequence matching using parallel brute force vs dynamic programming?*
- Recent JHU effort on integrating CUDA with SQL Server, using SQL UDF
- Galaxy correlation functions:  
400 trillion galaxy pairs!



# Tradeoffs Today

Stu Feldman: Extreme computing is about tradeoffs

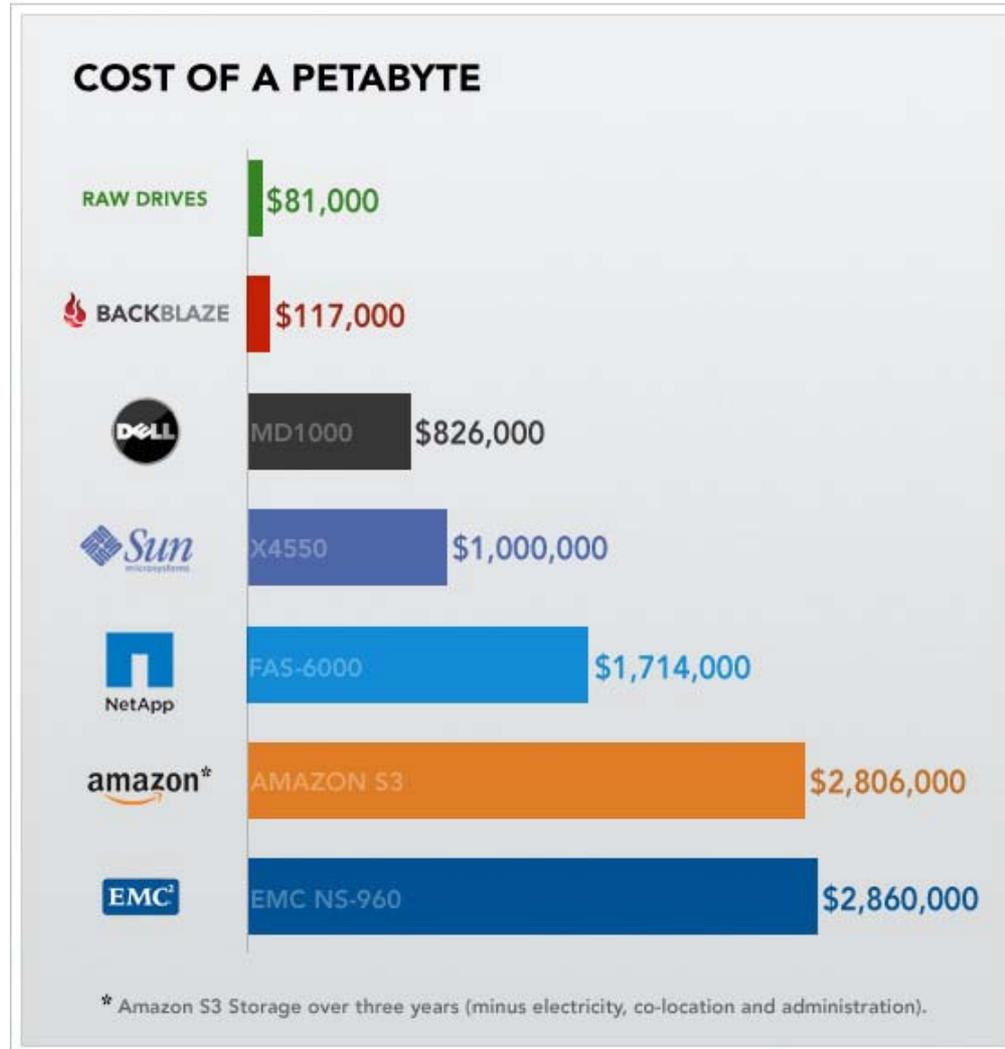
Ordered priorities for data-intensive science

1. *Total storage* (-> *low redundancy*)
2. *Cost* (-> *total cost vs price of raw disks*)
3. *Sequential IO* (-> *locally attached disks, fast ctrl*)
4. *Fast stream processing* (-> *GPUs inside server*)
5. *Low power* (-> *slow normal CPUs, lots of disks/mobo*)

The order will be different in a few years...and scalability may appear as well

# Cost of a Petabyte

From backblaze.com



# Petascale Computing at JHU

- Distributed SQL Server cluster/cloud w.
- 50 Dell servers, 1PB disk, 500 CPU
- Connected with 20 Gbit/sec Infiniband
- 10Gbit lambda uplink to UIC
- Funded by Moore Foundation, Microsoft and Pan-STARRS
- Dedicated to eScience, provide public access through services
- Linked to 1000 core compute cluster
- Room contains >100 of wireless temperature sensors



# GrayWulf Performance

- Demonstrated large scale scientific computations involving ~200TB of DB data
- DB speeds close to “speed of light” (72%)
- Scale-out over SQL Server cluster
- Aggregate I/O over 12 nodes
  - *17GB/s for raw IO, 12.5GB/s with SQL*
- **Scales to over 70GB/s for 46 nodes from \$700K**
- Cost efficient: \$10K/(GB/s)
- Excellent Amdahl number : 0.50
- But: we are already running out of space.....

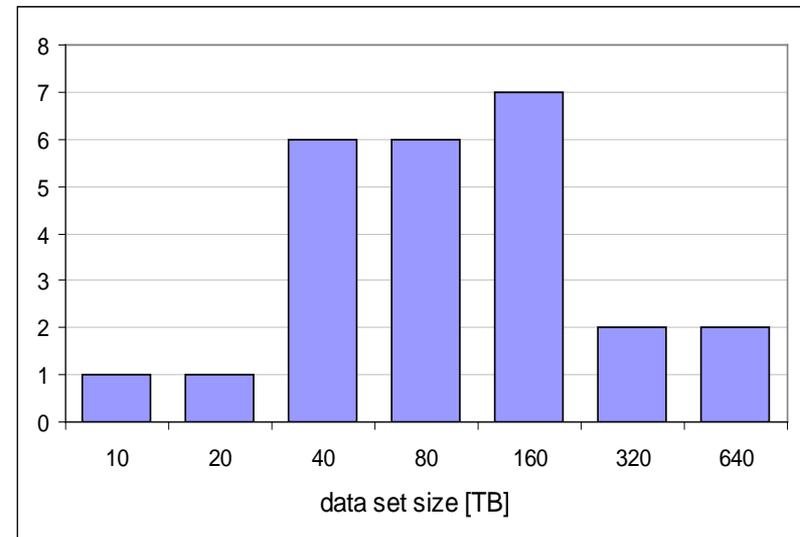
# Data-Scope

- Proposal to NSF MRI to build a new 'instrument' to look at data
- 102 servers for \$1M + about \$200K switches+racks
- Two-tier: performance (P) and storage (S)
- Large (5PB) + cheap + fast (460GBps), but ...
  - ..a special purpose instrument

	<i>1P</i>	<i>1S</i>	<i>90P</i>	<i>12S</i>	<i>Full</i>	
servers	1	1	90	12	102	
rack units	4	12	360	144	504	
capacity	24	252	2160	3024	5184	TB
price	8.5	22.8	766	274	1040	\$K
power	1	1.9	94	23	116	kW
GPU	6	0	540	0	540	TF
seq IO	4.6	3.8	414	45	459	GBps
netwk bw	10	20	900	240	1140	Gbps

# Proposed Projects at JHU

Discipline	data [TB]
Astrophysics	930
HEP/Material Sci.	394
CFD	425
BioInformatics	414
Environmental	660
Total	2823



19 projects total, data lifetimes between 3 mo and 3 yrs

# Short Term Trends

- Large data sets are here, solutions are not
  - *100TB is the current practical limit*
- National Infrastructure does not match power law
- No real data-intensive computing facilities available
  - *Some are becoming a “little less CPU heavy”*
- Even HPC projects choking on IO
- Cloud hosting currently very expensive
- Cloud computing tradeoffs different from science needs
- Scientists are “cheap”, also pushing the limit
  - *We are still building our own...*
  - *We will see campus level aggregation*
  - *May become the gateways to future cloud hubs*

# 5 Year Trend

- Sociology:
  - *Data collection in ever larger collaborations (VO)*
  - *Analysis decoupled, off archived data by smaller groups*
  - *Data sets cross over to multi-PB*
- Some form of a scalable Cloud solution inevitable
  - *Who will operate it, what business model, what scale?*
  - *How does the on/off ramp work?*
  - *Science needs different tradeoffs than eCommerce*
- Scientific data will never be fully co-located
  - *Geographic origin tied to experimental facilities*
  - *Streaming algorithms, data pipes for distributed workflows*
  - *“Data diffusion”?*
  - *Containernet (Church, Hamilton, Greenberg 2010)*

# Future: Cyberbricks?

- 36-node Amdahl cluster using 1200W total
- Zotac Atom/ION motherboards
  - *4GB of memory, N330 dual core Atom, 16 GPU cores*
- Aggregate disk space 43.6TB
  - *63 x 120GB SSD = 7.7 TB*
  - *27x 1TB Samsung F1 = 27.0 TB*
  - *18x.5TB Samsung M1= 9.0 TB*
- Blazing I/O Performance: 18GB/s
- Amdahl number = 1!
- Cost is less than \$30K
- Using the GPUs for data mining:
  - *6.4B multidimensional regressions in 5 minutes over 1.2TB*



# Summary

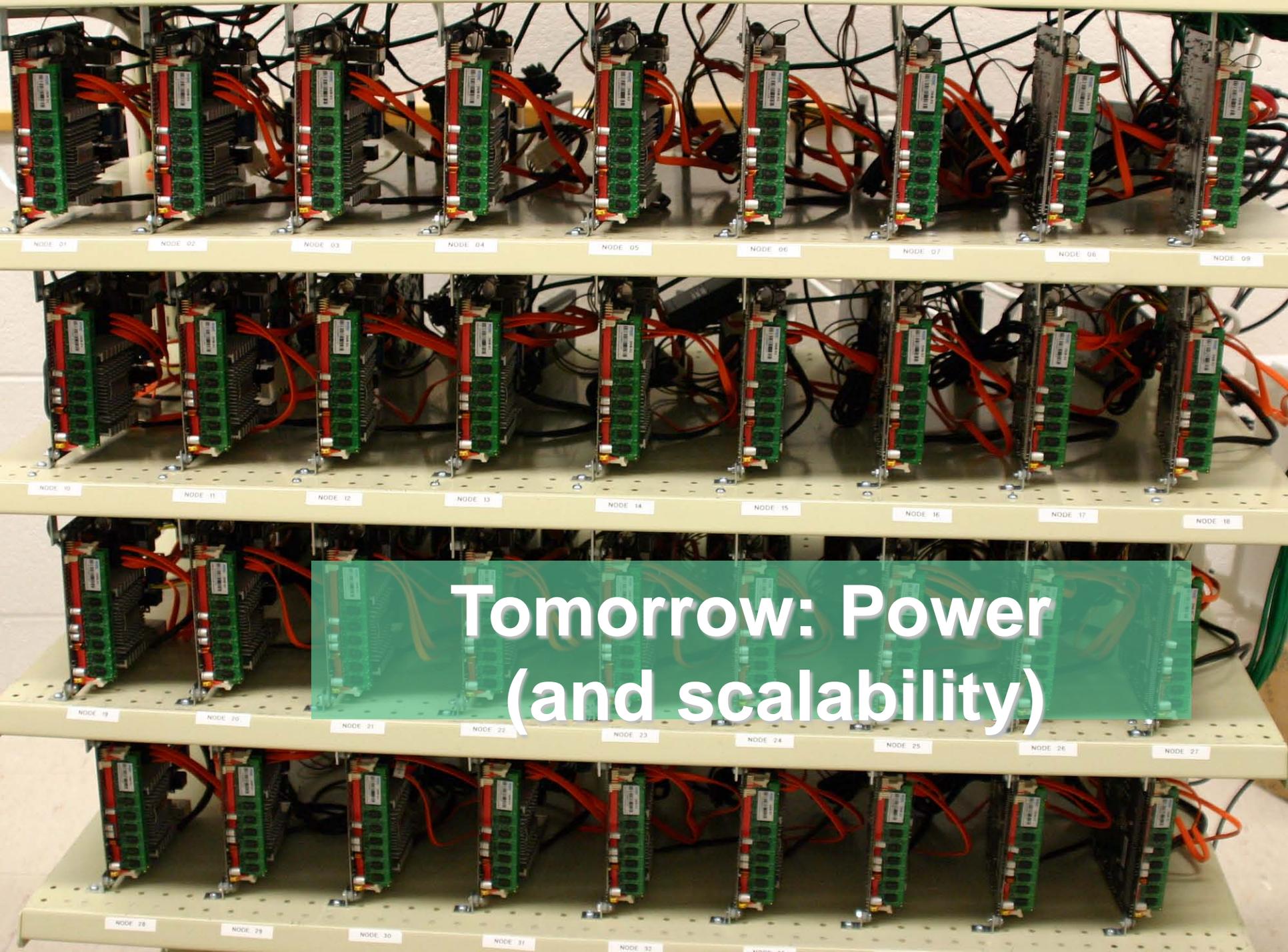
- Science community starving for storage and IO
  - *Data-intensive computations as close to data as possible*
- Real multi-PB solutions are needed NOW!
  - *We have to build it ourselves*
- Current architectures cannot scale much further
  - *Need to get off the curve leading to power wall*
  - *Multicores/GPGPUs + SSDs are a disruptive change!*
- Need an objective metrics for DISC systems
  - *Amdahl number appears to be good match to apps*
- Future in low-power, fault-tolerant architectures
  - *We propose scaled-out “Amdahl Data Clouds”*
- A new, Fourth Paradigm of science is emerging
  - *Many common patterns across all scientific disciplines*

**Yesterday: CPU cycles**

A 3D rendered scene of a server room. In the foreground, there are several server racks. The racks are dark grey or black with lighter grey panels. The floor is made of light-colored square tiles. The background shows more server racks receding into the distance. The text 'Yesterday: CPU cycles' is written in a bold, white, sans-serif font across the middle of the image.

# Today: Data Access





# Tomorrow: Power (and scalability)