

# The Big Stuff

# Focus on Big

- Lots of talks concerning or in anticipation of lots of data and/or lots of computation
- No surprise given session names: “Large Observatory Challenges”, “Grid and Grid Virtualization”, “Cross Catalog Matching”,...
- The expected buzzwords: Grid, Cloud, VO....
- And (for me) some new ones: Hadoop, Pig, GPU

# Typical “Big” Talk



## Amdahl's Laws and Extreme Data-Intensive Scientific Computing

Alex Szalay  
The Johns Hopkins University

### Summary

- Large data sets are here, solutions are not
  - *100TB is the current practical limit*
- Science community starving for storage and IO
- No real data-intensive computing facilities available
  - *Changing with Dash, Gordon, Data-Scope, GrayWulf...*
- Even HPC projects choking on IO
- Real multi-PB solutions are needed NOW!
- Cloud hosting currently very expensive
- Cloud computing tradeoffs different from science needs
- Scientists are “frugal”, also pushing the limit
- Current architectures cannot scale much further
- Astronomy representative for science data challenges

# Cloud



## On-line Access and Visualization of Multi- dimensional FITS Data

Pavol Federl

Institute for Space Imaging Science  
University of Calgary



## Overview

- CyberSKA project
  - Develop cyberinfrastructure for SKA
  - Online accessible tools
  - HW/SW requirements: computer + modern browser
- Cyber SKA web portal  
[www.cyberska.org](http://www.cyberska.org)

# Cloud cont'd

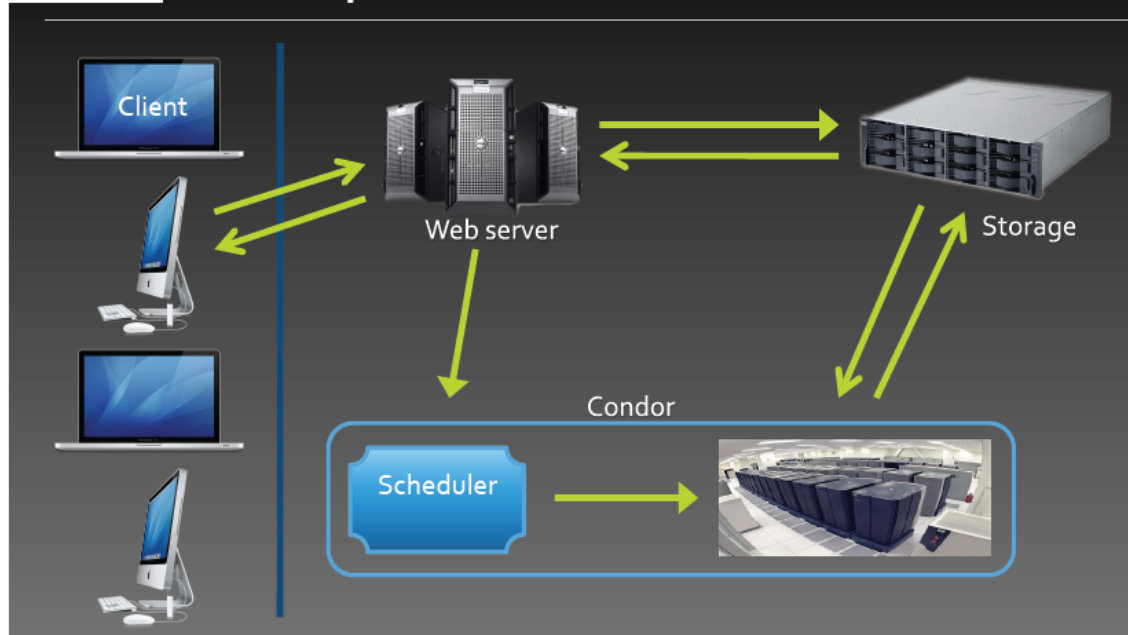


## Implementation Details

- Client side
  - HTML 5
  - JavaScript
  - Ajax
- Server side
  - PHP
  - Condor Pool
  - C++
  - NFS



## Implementation Details



# Grid

## Measure transverse motion of 730.000 stars - 1



- How many stars are in the Pleiades besides the famous ones?



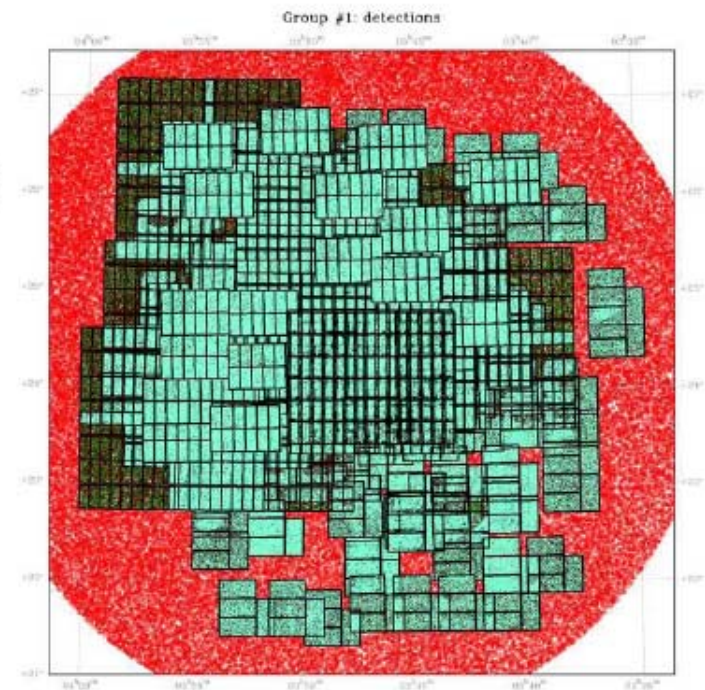
- Difficult to tell just from the right image
- The only good way to find out is to look for common transverse motion

# Grid cont'd

## Measure transverse motion of 730.000 stars - 2



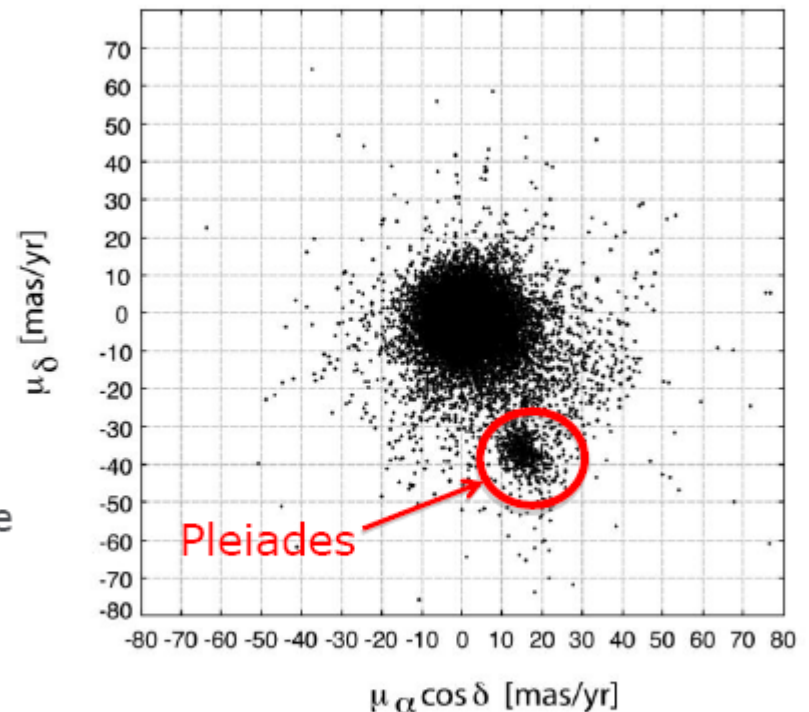
- Use of Scigrad to process 520Gb of wide field images obtained from various space and ground based telescopes (Subaru, HST, CFHT, Isaac Newton Telescope) over a 12 years period
- The multi-epoch images were used to derive the transverse motion of every star ( $\sim 730.000$  in total) present in the 6deg x 6deg field of view of our observations
- The usage of Scigrad was critical, as it requires:
  - vast amount of storage
  - fast multi-threaded computers to extract the source photometry and astrometry
  - fast multi-threaded computers and vast amounts of RAM to cross-match
  - the multi-epoch catalogs and derive the kinematics of each star



H.Bouy et al

# It worked!

- The result is a vector point diagram of the motion of all 730.000 star located in the Pleiades cluster. The diagram shows the Field and background objects are distributed randomly around (0,0) mas/yr, while the Pleiades members are all co-moving and form the locus near (15,-35)
- This allowed to unambiguously identify several thousands of members (when less than 1000 were known to date) down to the planetary mass regime.
- The scientific outcome will be extremely valuable and rich, from refining the mass function of the cluster, identifying planetary mass objects, and detailed studies of the internal dynamics.
- This new technique could be repeated for other regions of the sky



motion in RA (x-axis) and Dec (y-axis)  
H.Bouy et al



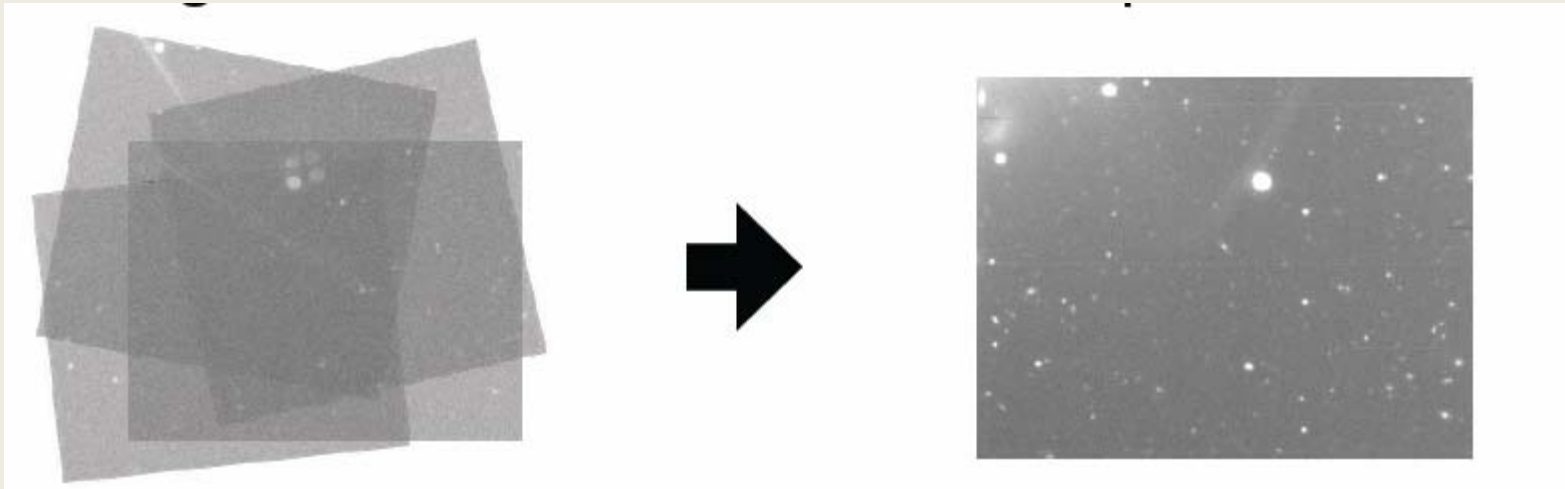
# Hadoop and Pig

- Hadoop is an Apache open-source project used by, e.g., Yahoo!, Facebook, Amazon.
- “The Apache Hadoop project develops open-source software for reliable, scalable, distributed computing.”
- Main components include
  - HDFS: distributed file system
  - MapReduce: software framework for distrib. processing
- Intellectual basis apparently from Google. Implemented in Java.



# Hadoop and Pig

*Astronomical Imaging with Hadoop* described how a Hadoop installation was use to do image addition:



Given region and images covering it  
crop, background-subtract, apply PSF, weight, etc.

# Pig

- High-level language wrapping Hadoop, developed by Yahoo!
- See *Pig as a Solution for Accessing Peta-scale Astronomical Datasets* for more background on Pig and their db access problem.



# Why Pig?

Because I bet you can read the following Pig script

```
top_5.pig
users = load 'users.csv' as (username: chararray, age: int);
users_1825 = filter users by age >= 18 and age <= 25;
pages = load 'pages.csv' as (username: chararray, url: chararray);
joined = join users_1825 by username, pages by username;
grouped = group joined by url;
summed = foreach grouped generate group as url, COUNT(joined) AS views;
sorted = order summed by views desc;
top_5 = limit sorted 5;
store top_5 into 'top_5_sites.csv';
```



# Why Pig?

## The same in Hadoop MapReduce

```
Mapper {
  @Override
  public void setup(Context context) {
    // ...
  }
  @Override
  public void map(LongWritable key, Text value, Context context) {
    // ...
  }
  @Override
  public void reduce(LongWritable key, List values, Context context) {
    // ...
  }
}
```

```
Mapper {
  // ...
  public void map(LongWritable key, Text value, Context context) {
    // ...
  }
  public void reduce(LongWritable key, List values, Context context) {
    // ...
  }
}
```

```
Mapper {
  // ...
  public void map(LongWritable key, Text value, Context context) {
    // ...
  }
  public void reduce(LongWritable key, List<Text> values, Context context) {
    // ...
  }
}
```

# GPU

- Graphics Processing Unit
- Commodity item. They're improving faster than CPUs (gaming apps a strong incentive)
- Very, very parallel
- Can program like general-purpose cpu
- Only single-precision, but that may change soon. Single is often good enough (so people say).

# GPU cont'd

- At least 4 talks (out of ~25 total) entirely on GPU apps:
  - *Massively Parallel Fourier-Space Cross-Correlation for Analyzing Highly Dimensional Time Series Databases*
  - *Fitting Galaxies on GPUs*
  - *Distributed GPU Volume Rendering of ASKAP Spectral Data Cubes*
  - And my favorite, *Astrophysical N-body Simulation on a Cluster of GPUs*

# Challenges

- To get a high efficiency on GPUs for hierarchical  $O(N \log N)$  or  $O(N)$  method (not on brute  $O(N^2)$  method)

Treecode, FMM

- using large amount of GPUs  
414,720 = 240 \* 3 \* 576 FP units

Multiple walk

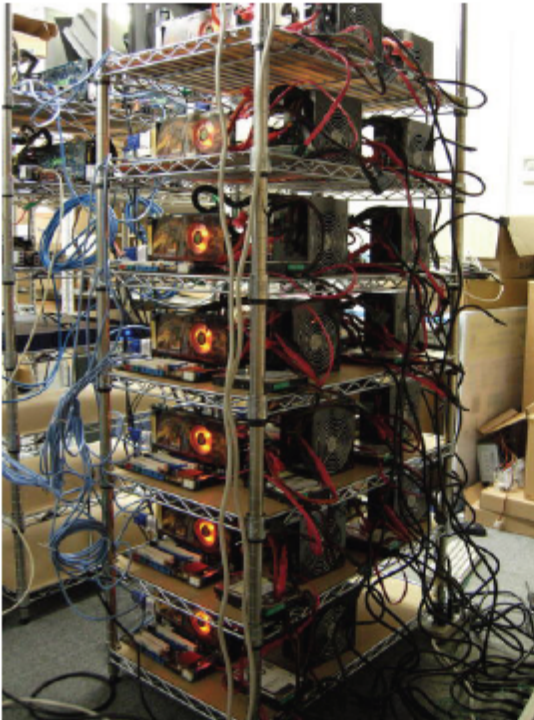
- To get a good scalability on commodity network (GbE)

Delegated Alltoallv



# History of our GPU cluster

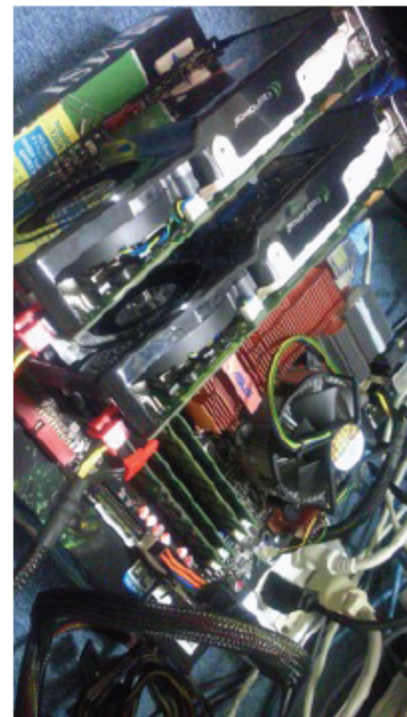
## Mar 2008



Host: Core2Quad 2.4 GHz .... x 32

GPU: GeForce 8800GT .... x 32

# Nov 2008



Host: Core2Quad 2.4 GHz .... x 128  
GPU: GeForce 8800GTS .... x 256

~ 40 Tflops in cosmology sim.

# Aug 2009



Power supply 600A → 2000A

Host: Core2Quad 2.4 GHz .... x 166

GPU: GeForce 9800GTX+ .... x 256

GPU: GeForce GTX295 .... x 33



出島

DEGIMA

cluster

J. Bogart

2010年11月11日 木曜日

# References

<http://adass2010.cfa.harvard.edu/ADASS2010/Program/presentations.html>