# BASIN

## *B*eowulf *A*nalysis *S*ymbolic *IN*terface

### *A multi-user environment for parallel data analysis and visualization*

Enrico Vesperini

Doug Jones

Dave Goldberg

Steve McMillan (PI)

and the BASIN team
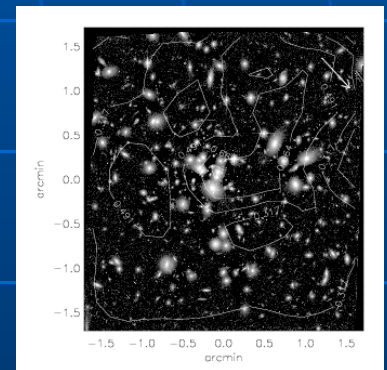
*Department of Physics*

*Drexel University*

SciPy2007

August 16,2007

# Who we are…

Theoretical and observational astrophysicists working with large ($10^2$-$10^3$ GB) datasets from simulations and observations:

- stellar dynamics, N-body simulations of dense stellar systems, globular star clusters;
- gravitational lensing, dark matter distribution
- observational extra-galactic astronomy, Sloan Digital Sky Survey-SDSS.



Globular Cluster NGC 6093

# What we use:
## Beowulf Clusters

- Beowulf clusters: attractive (low-cost) parallel systems

- Several packages exist for serial data analysis (e.g., in observational astronomy, IRAF, MIDAS).

- For parallel data analysis each group tends to (re-) develop its own set of tools:
  - Expertise in parallel tools/algorithms.
  - Large time investment.

# What is it?

- Integrated suite of tools for parallel data analysis and visualization.

- Multi-user environment for interactive data analysis and visualization.

# Why?

- Easy and transparent parallel data analysis.
- Avoid redundant development of functions commonly used in data analysis.

# Where?

## BASIN is freely (GPL) available at
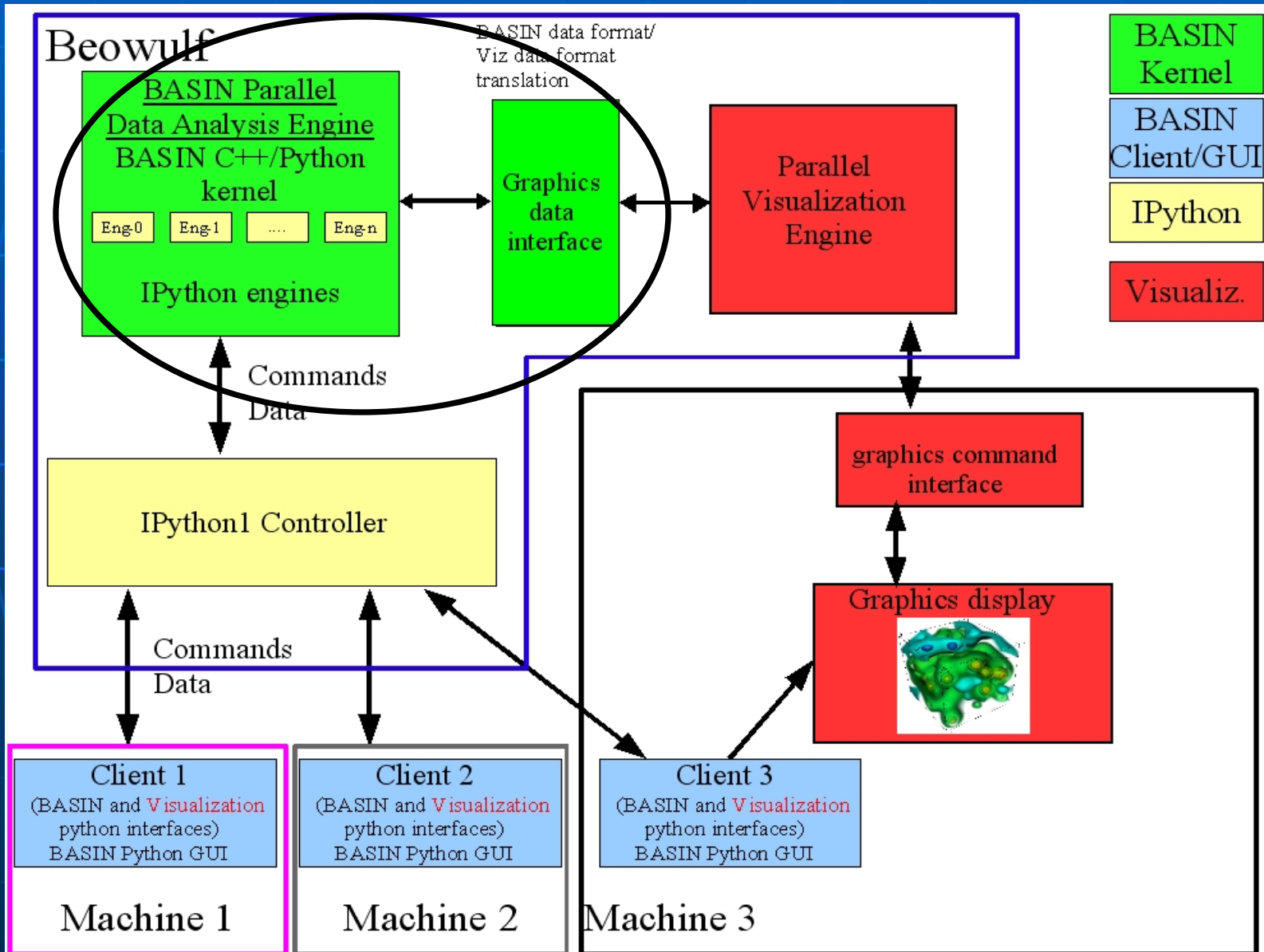
## http://www.physics.drexel.edu/BASIN

# BASIN DEMO

1) Start the computational engine on a remote parallel server.

2) Connect a local client to the computational engine.

3) Read file distribute data.

4) Parallel calculation of new attributes.

5) Visualize distributed data (using VisIt-LLNL)

6) Transfer data on the local machine and plot with one of the standard Python plotting packages (Matplotlib).
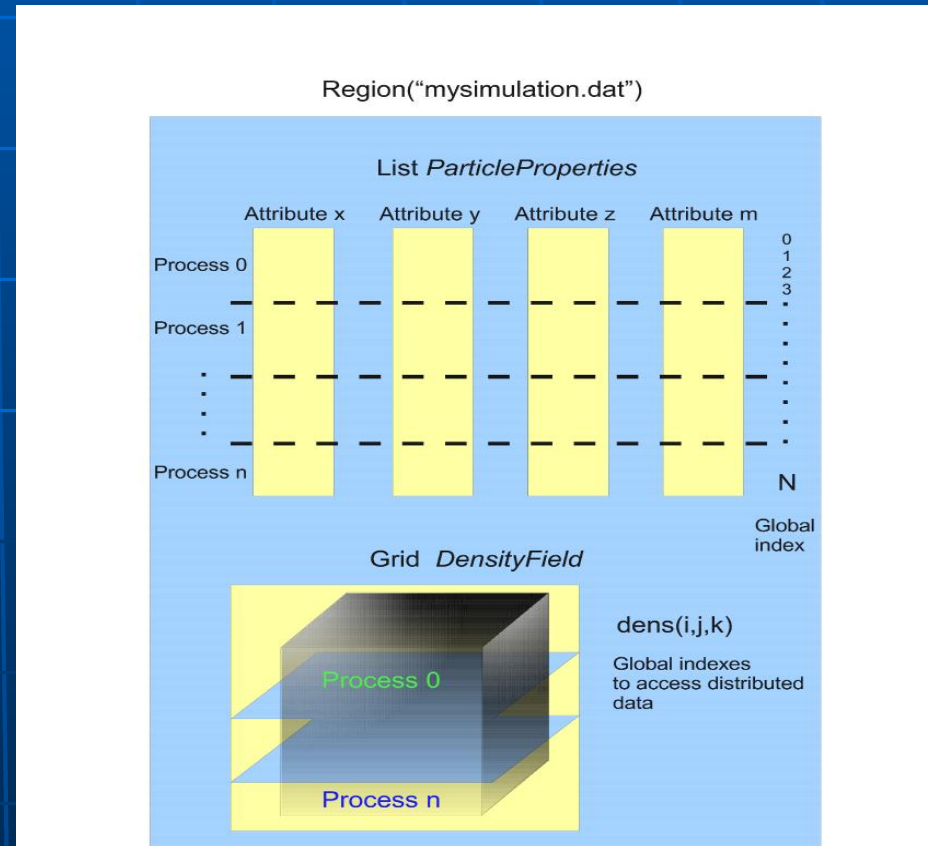
# BASIN architecture

# BASIN architecture

# Data Analysis Engine

## *BASIN kernel : classes and functions for data distribution and parallel data operations.*
### *(C++/MPI)*

Objects visible to the user:

- ***Region*** defined by the data file read

- ***Data*** :
  - ***Grid***
  - ***List***

- ***Attribute***

# Data distribution and parallel operations
## Attribute class

- Hide the complexity of data distribution, retrieval and parallel operations.
- Tools to ease the parallel data distribution and management on a distributed memory machine.
- Shared-memory view of data in a distributed memory machine.
- Arbitrary n-dimensional arrays of primitive datatypes or user-defined structures and objects.
- Parallel math ( elementwise and reduction) and logical operations on distributed data.
- Each process can locate its own block of data.

# Data distribution and parallel operations
## Attribute class

- Hide the complexity of data distribution, retrieval and parallel operations.
- Tools to ease the parallel data distribution and management on a distributed memory machine.
- Shared-memory view of data in a distributed memory machine.
- Arbitrary n-dimensional arrays of primitive datatypes or user-defined structures and objects.
- Parallel math ( elementwise and reduction) and logical operations on distributed data.
- Each process can locate its own block of data.

# Data distribution and parallel operations
## Attribute class

- Hide the complexity of data distribution, retrieval and parallel operations.
- Tools to ease the parallel data distribution and management on a distributed memory machine.
- Shared-memory view of data in a distributed memory machine.
- Arbitrary n-dimensional arrays of primitive datatypes or user-defined structures and objects.
- Parallel math ( elementwise and reduction) and logical operations on distributed data.
- Each process can locate its own block of data.

# Data distribution and parallel operations
## Attribute class

- Hide the complexity of data distribution, retrieval and parallel operations.

- Tools to ease the parallel data distribution and management on a distributed memory machine.

- Shared-memory view of data in a distributed memory machine.

- Arbitrary n-dimensional arrays of primitive datatypes or user-defined structures and objects.

- Parallel math ( elementwise and reduction) and logical operations on distributed data.

- Each process can locate its own block of data.

# Data distribution and parallel operations
## Attribute class

- Hide the complexity of data distribution, retrieval and parallel operations.
- Tools to ease the parallel data distribution and management on a distributed memory machine.
- Shared-memory view of data in a distributed memory machine.
- Arbitrary n-dimensional arrays of primitive datatypes or user-defined structures and objects.
- Parallel math ( elementwise and reduction) and logical operations on distributed data.
- Each process can locate its own block of data.

# Data distribution and parallel operations
## Attribute class

- Hide the complexity of data distribution, retrieval and parallel operations.
- Tools to ease the parallel data distribution and management on a distributed memory machine.
- Shared-memory view of data in a distributed memory machine.
- Arbitrary n-dimensional arrays of primitive datatypes or user-defined structures and objects.
- Parallel math ( elementwise and reduction) and logical operations on distributed data.
- Each process can locate its own block of data.

# Data distribution and parallel operations
## Attribute class

```
reg0 = Region ("/home/vesperin/starclust.dat")
m = reg0.get_attribute ("mass")
x = reg0.get_attribute ("x")
y = reg0.get_attribute ("y")
z = reg0.get_attribute ("z")
print m[10]
m[120]=12.3
total_Mass=sum(m)
logm=log10(m)
Hm=where(m>5,m,0)
logr=log10(sqrt(x*x+y*y+z*z))
```

# Center of Mass

- ## C/ MPI
- ## BASIN(C++/Python)

```
……m, x…..

MPI_Scatter(&m[0], n_loc, MPI_DOUBLE,
&m_loc[0],n_loc,MPI_DOUBLE, 0,
MPI_COMM_WORLD);
MPI_Scatter(&m[0], n_loc, MPI_DOUBLE,
&m_loc[0],n_loc,MPI_DOUBLE, 0,
MPI_COMM_WORLD);


for(int i = 0; i < n_loc; i++)
{
x_sum_loc += x_loc[i] * m_loc[i];
m_sum_loc += m_loc[i];
}

MPI_Reduce(&x_sum_loc, &x_sum, 1,
MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
MPI_Reduce(&m_sum_loc, &m_sum, 1,
MPI_DOUBLE, MPI_SUM, 0,MPI_COMM_WORLD);

    if(myRank == 0) {
      com[0] = x_sum / m_sum;
      }
```

```
com=sum(m*x)/sum(m)
```

# Data Analysis Engine

## BASIN kernel : Scientific Packages

- *Cosmology*
- *Stellar dynamics*
- *Statistics*
- *FFT* (FFTW, http://www.fftw.org/)
- *Coordinate transformations*

# BASIN architecture

# BASIN Python User Interface and GUI

- BASIN Python interface created with Boost Python.

- A remote Python client can invoke BASIN commands to be executed by the Data Analysis Engine.

- Multiple distributed clients can connect to the same BASIN Data Analysis Engine and share the same data ( *IPython/IPython1* F.Perez, B.Granger `http://ipython.scipy.org/moin/IPython1`)

# BASIN Python User Interface and GUI



- Connection to the remote server

- I/O

- Save/import session history

# BASIN Python User Interface and GUI

# BASIN Python User Interface and GUI



• Visualization packages

# BASIN Python User Interface and GUI

# BASIN Python User Interface and GUI



•Scientific packages

# BASIN Python User Interface and GUI



BASIN Session data structure

Reference to `Attribute` and "free" `Attributes`

Python window

Python output

Server/local client switch

# BASIN architecture

# Visualization

## VisIt

developed at LLNL
http://www.llnl.gov/visit



- Visualization of large distributed datasets (structured and unstructured meshes)

- Parallel visualization engine

- Available in BASIN (in collaboration with Brad Whitlock at LLNL)

# Visualization

## Partiview

Developed at NCSA
http://virdir.ncsa.uiuc.edu/partiview/



•Visualization of 4D particle datasets

• BASIN/Partiview interface: work in progress ( in collaboration with Stuart Levy, Matt Hall  at NCSA)

# Visualization

## Gnuplot



- Simple plotting package based on the Gnuplot API.

- Only for development purposes.

- Available in BASIN

# Visualization

Data transferred to the client machine: all the Python plotting packages (e.g. Matplotlib, Gnuplot, Chaco,  etc.)

# Summary

*Goals:*

Ease access to parallel data analysis

Avoid redundant development

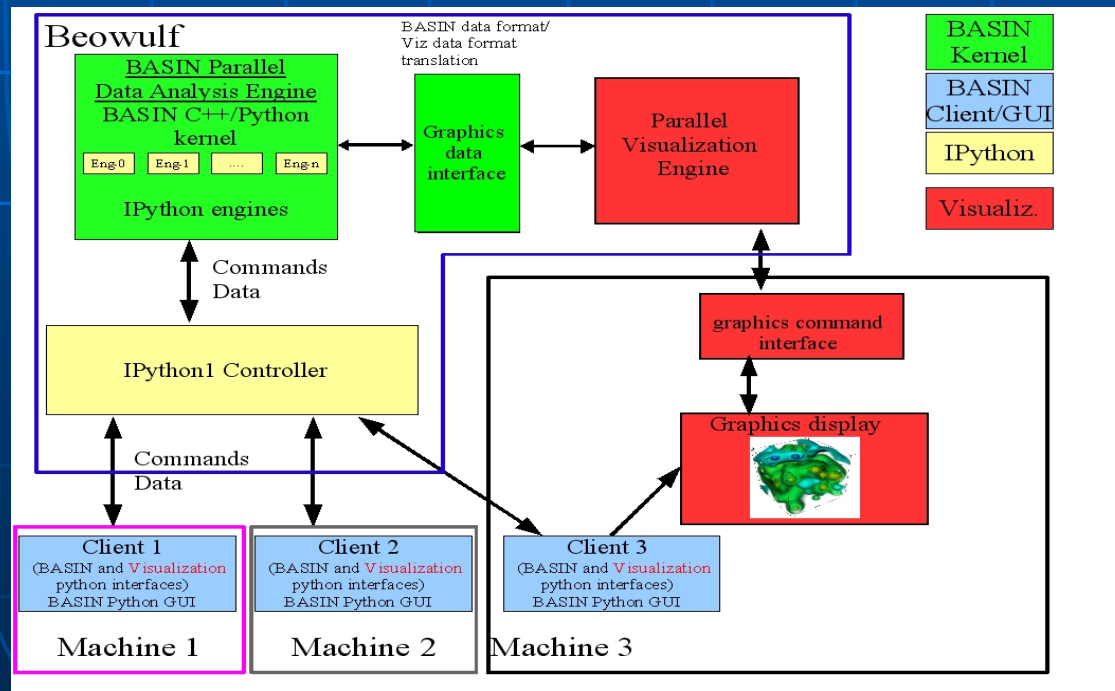Interactive and multi-user parallel data analysis

# Summary

*What we have:*

Kernel for parallel data management and operations (C++/Python)

Scientific packages

Interface to a few existing visualization packages

# Summary

*What to look for next:*

Increase science scope beyond astrophysics

Extend visualization options

Two-way communication with visualization packages

Improve ease of use and installation