**Thorin Tabor**
JupyterCon 2017
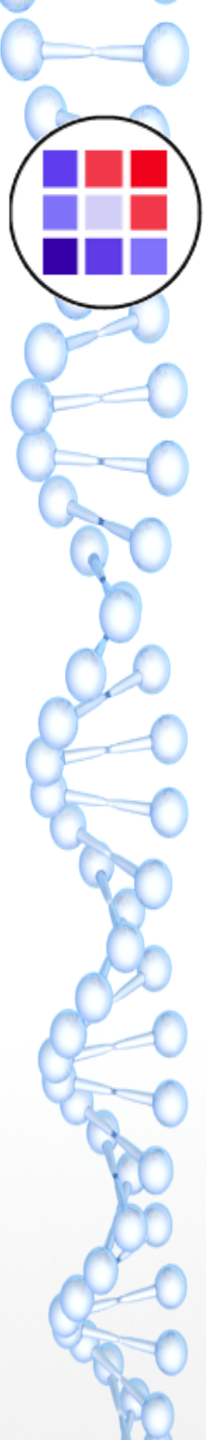
# GenePattern Notebooks
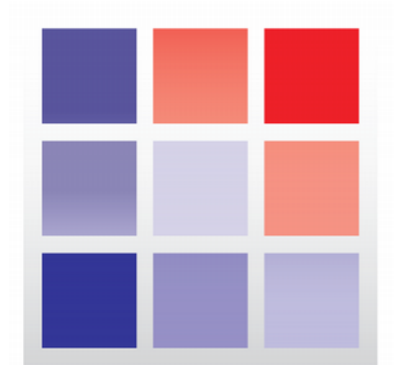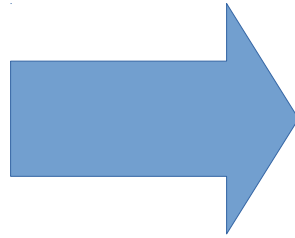
Jupyter for Integrative Genomics

UC San Diego

BROAD INSTITUTE

# From Jupyter to GenePattern

# GenePattern Notebook

# Two Open Source Projects

# Platform for Reproducible Bioinformatic Research

- First public release in 2004 (similar footing to IPython)

- Open Source

- ~50,000 registered users

- Public server runs ~4,000 analyses per week

- Community-contributed methods

  - CRISPR analysis

  - Bisulfite sequencing

  - Flow cytometry

  - RNAi screens



**genepattern.org**



**genepattern.broadinstitute.org**



**gparc.org**

# Analysis Tool Repository

| | | | |
|---|---|---|---|
| Copy Number Divide by Normals | GSEA | Variation Filter | Cuffdiff |
| GISTIC | CBS | k-Nearest Neighbors | MutSigCV |
| Classification and Regression Trees | Support Vector Machines | Hierarchical Clustering | Picard Sort Sam |
| TopHat | Expression File Creator | Metagene Projection | RNASeQC |

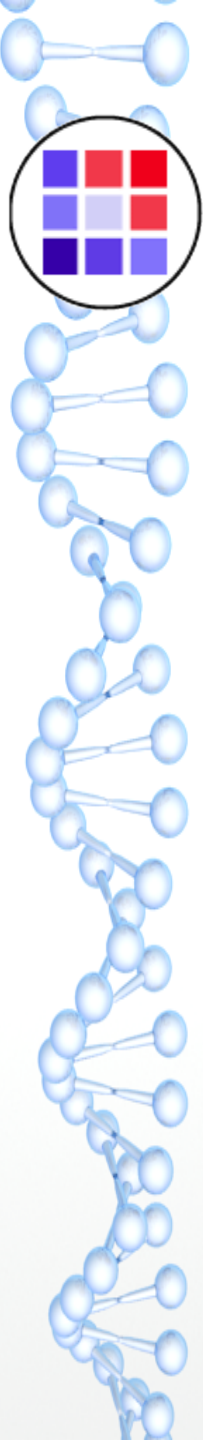# Custom Modules & Pipelines

## Modules

### Hierarchical Clustering Files

HCL.jar
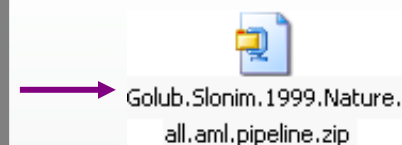cluster.sh
ant.jar
gp-modules.jar
Jama-1.0.2.jar

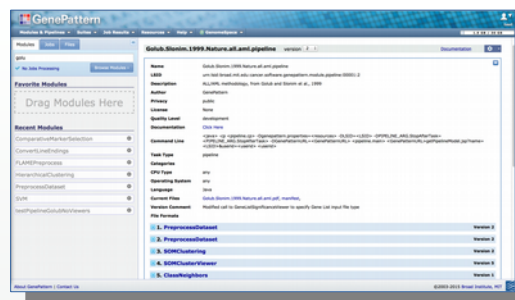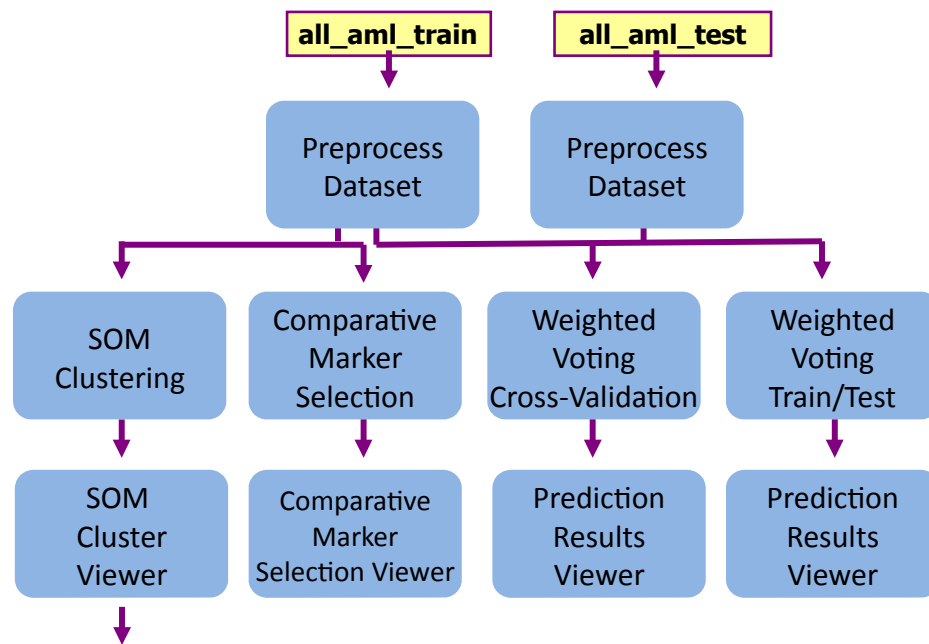### Documentation

HierarchicalClustering.pdf

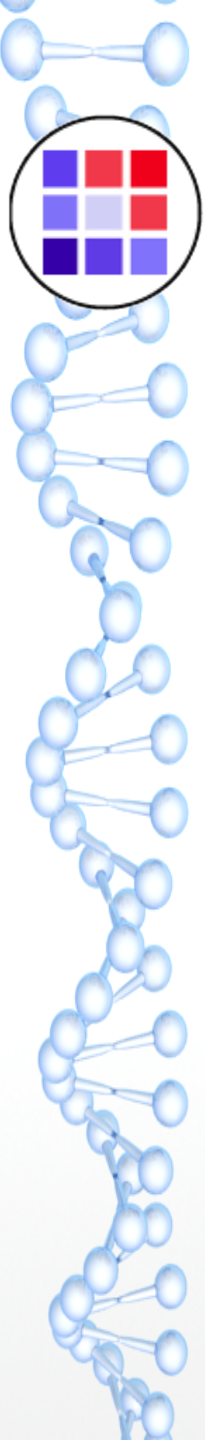### Parameter descriptions

-f <input.filename>
<log.transform>
<row.center>
<row.normalize>
<column.center>
<column.normalize>
-u <output.base.name>
-e <column.distance.measure>
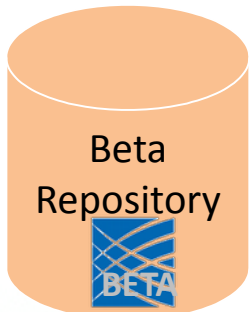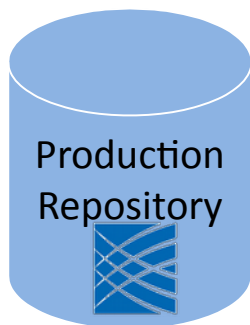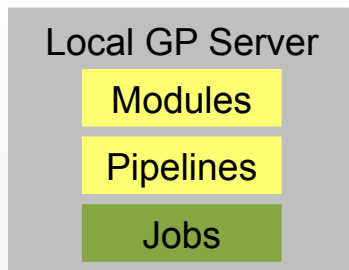-g <row.distance.measure>
-m <clustering.method>

## Pipelines

# Web Server Architecture

# Programmatic APIs

- Libraries for Python, R, MATLAB & Java

- REST API

- Used to back portals and other web applications

# User-Friendly Interface

- Permits complex analyses without the need for a coding background

# GenePattern Notebook Jupyter Extensions

# Complete Research Narrative

- Leverage the best of Jupyter and GenePattern

- Interleave text, visualizations, graphics and analytical aspects

# SVM Example #1

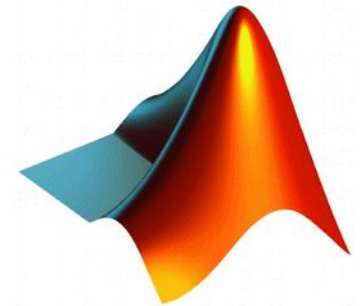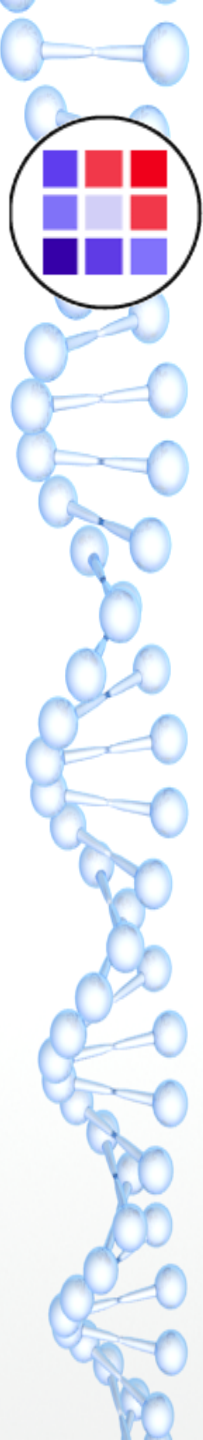```python
%matplotlib inline

import numpy as np
import matplotlib.pyplot as plt
from sklearn import svm
```

```python
# Load the training data
train_data = None
with open('/home/thorin/datasets/all_aml_train.gct', 'r') as td:
    raw_txt = td.read()
    train_data = np.genfromtxt(fname=raw_txt, delimiter='\t', dtype=None, comments=None)
```

```python
# Load the training classes
train_classes = None
with open('/home/thorin/datasets/all_aml_train.cls', 'r') as tc:
    raw_txt = tc.read()
    train_classes = np.genfromtxt(fname=raw_txt, delimiter=' ', dtype=None, comments=None)
```

```python
# Slice the data for SVM fitting
X = train_data.data[:, :2]
y = train_classes
```

```python
# Create an instance of SVM and fit out data. Do not scale the data.
C = 1.0 # SVM regularization parameter
svc = svm.SVC(kernel='linear', C=1,gamma='auto').fit(X, y)
```

```python
# Create a mesh to plot in
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
h = (x_max / x_min)/100
xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
 np.arange(y_min, y_max, h))
```

```python
plt.subplot(1, 1, 1)
```

# SVM Example #2

# GenePattern Cells

Auth Cell

Analysis Cell

Job Cell

# Authentication Cells



**GenePattern** *Login*

**GenePattern Server**

Broad Institute

**GenePattern Username**

Username

**GenePattern Password**

Password

Log into GenePattern    Register an Account



**GenePattern** *tabor*    http://genepattern.broadinstitute.org/gp

-- Sun 5:00 pm -- Update: The job queue is back online and accepting new jobs. For best results you should cancel any jobs which you had started before today at 5:00 pm. We can not make any guarantees about results obtained for jobs that had not yet completed before the start of the maintenace window.Thanks,The GenePattern Team -- Sat 5:00 pm -- Update: The job queue is not yet ready to accept new jobs. Please refrain from starting new jobs until further notice. We expect it to be ready during the day Sunday.Thanks,The GenePattern Team Important message: The GenePattern Server will go offline for quarterly maintenance just before 8:00 am, Saturday March 5. We expect the maintenance to last the majority of the day.Thanks,The GenePattern Team -- March 7 -- New Blog Post: Older Java Applet Visualizers Blocked by Default in Updated FirefoxOlder Java Applet visualizers are no longer supported in Chrome. Please read our blog post for more information.

Experiencing a bug? Have thoughts on how to make GenePattern Notebook better? Let us know by leaving feedback.    Leave Feedback

# Analysis Cells

# Job Cells

# Python Function GUI



- Turn any Python function into an interactive user interface

# Rich Text Markdown Editor



- No markup knowledge required
- Generates HTML / markdown
- Available as a separate extension

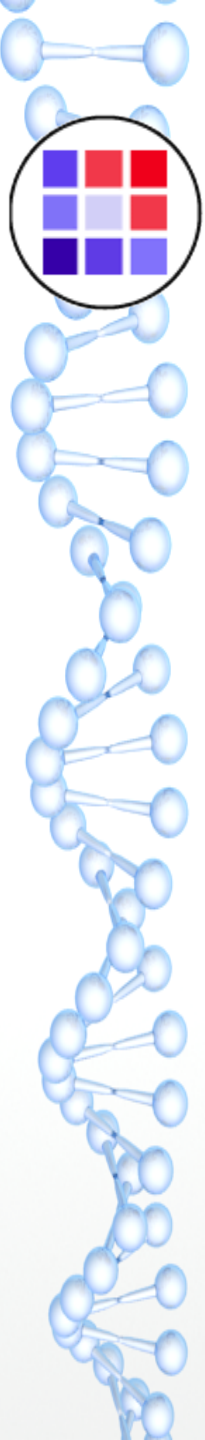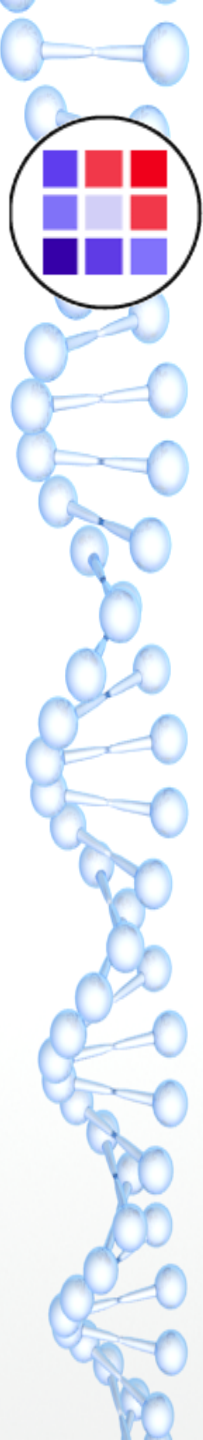# Notebook Tool Manager

# Behind the Scenes

- Interactive widgets use the Jupyter widget framework (ipywidgets, traitlets)

- Can use Python variables as input

- Not limited by GenePattern analyses

# GenePattern Python Library

- Complete programmatic access

- Automatic integration with GenePattern cell data

```
import gp

# Create a GenePattern server proxy instance
gpserver = gp.GPServer('http://localhost:8080/gp','myusername', 'mypassword')

# Obtain GPTask by module name
module = gp.GPTask(gpserver, "PreprocessDataset")

# Load module parameter data
module.param_load()

# Create a job specification
job_spec = module.make_job_spec()

# Upload a file to the server
uploaded_file = gpserver.upload_file("file_name", "/path/to/the/file/file_name")
job_spec.set_parameter("input.filename", uploaded_file.get_url())

# Submit the job to the GenePattern server
job = gpserver.run_job(job_spec)
```
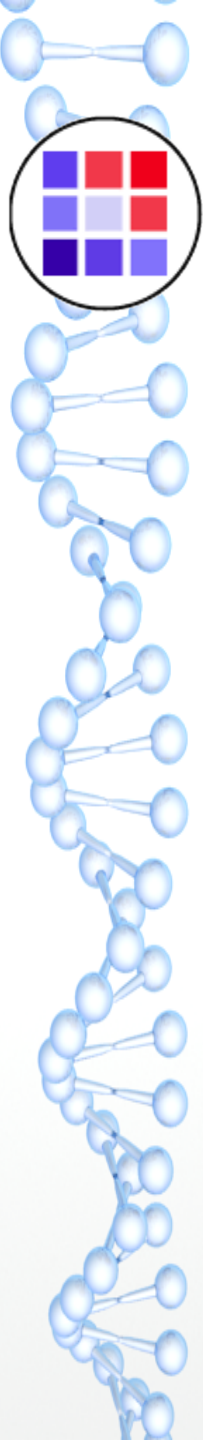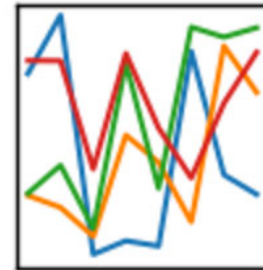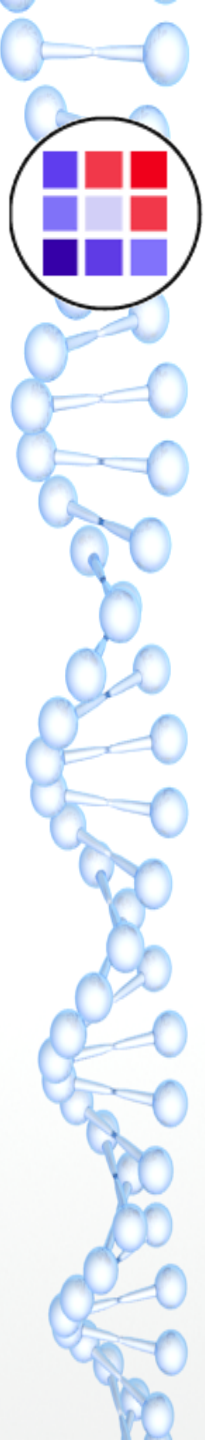
# GenePattern Data Tools

- Easily import common bioinformatic data formats as pandas DataFrames

- Work with GenePattern files using popular Python libraries



pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

# GenePattern Notebook Repository



http://genepattern-notebook.org

# Publish & Share Notebooks



- Publish notebooks to the GenePattern Notebook Repository.
- Browse available notebooks.

# Installing the Extension

- ## PyPI

  - pip install genepattern-notebook

- ## Anaconda Cloud

  - conda install -c genepattern genepattern-notebook

- ## DockerHub

  - docker pull genepattern/genepattern-notebook

# Jupyter Ecosystem

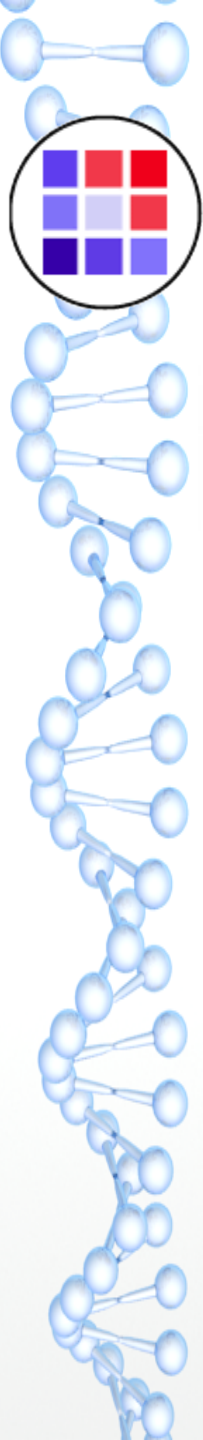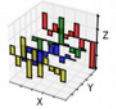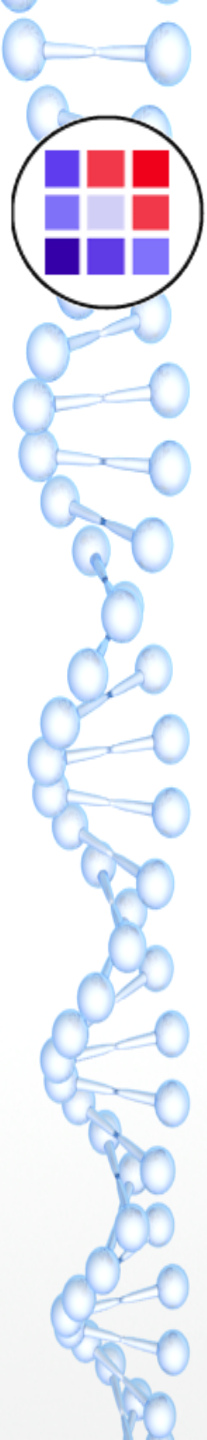# Acknowledgments

## GenePattern Team

Peter Carr

David Eby

Barbara Hill

Edwin Juarez

Ted Liefeld

Michael Reich

Thorin Tabor

Helga Thorvaldsdottir

## PI

Jill Mesirov

# Resources

GenePattern Notebook
genepattern-notebook.org

GenePattern
genepattern.org

Public GenePattern server
genepattern.broadinstitute.org

Indiana University GenePattern server
gp.indiana.edu

GenePattern Archive (GPArc)
gparc.org

GenePattern Twitter
@genepattern

GenePattern GitHub
github.com/genepattern

GenePattern DockerHub
hub.docker.com/r/genepattern