IP[y]:
IPython

jupyter

# **Reproducible science with Jupyter**

Changing our publication models
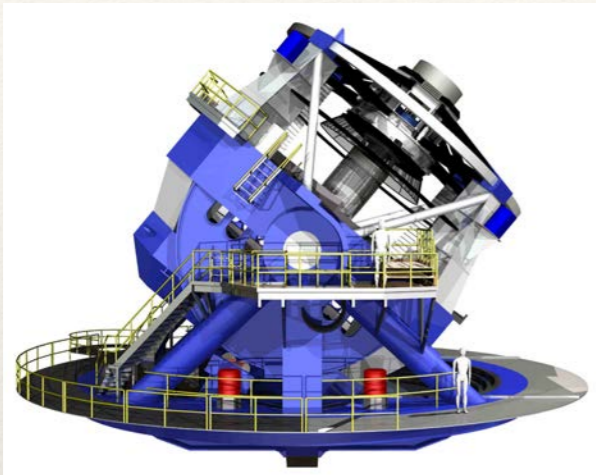
Fernando Pérez
(@fperez_org & fperez@lbl.gov)

LBL & UC Berkeley

BERKELEY LAB

Berkeley
UNIVERSITY OF CALIFORNIA

BIDS
BERKELEY INSTITUTE
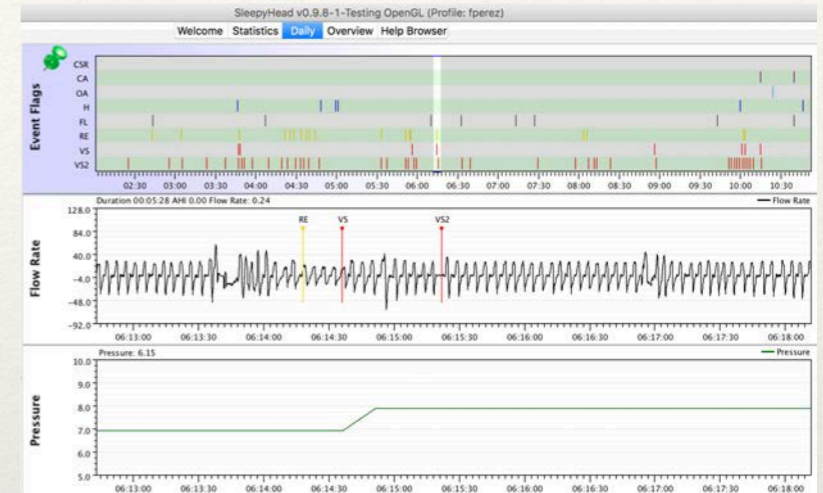FOR DATA SCIENCE

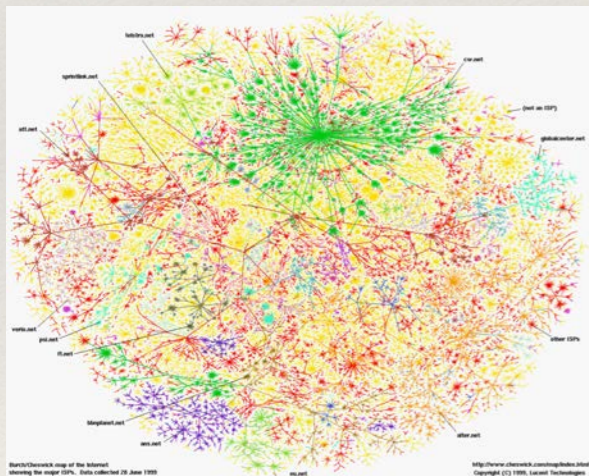# Every research discipline is now awash in data



Astronomy: LSST



Physics: LHC



Personalized, data-driven medicine



Sociology: The Web



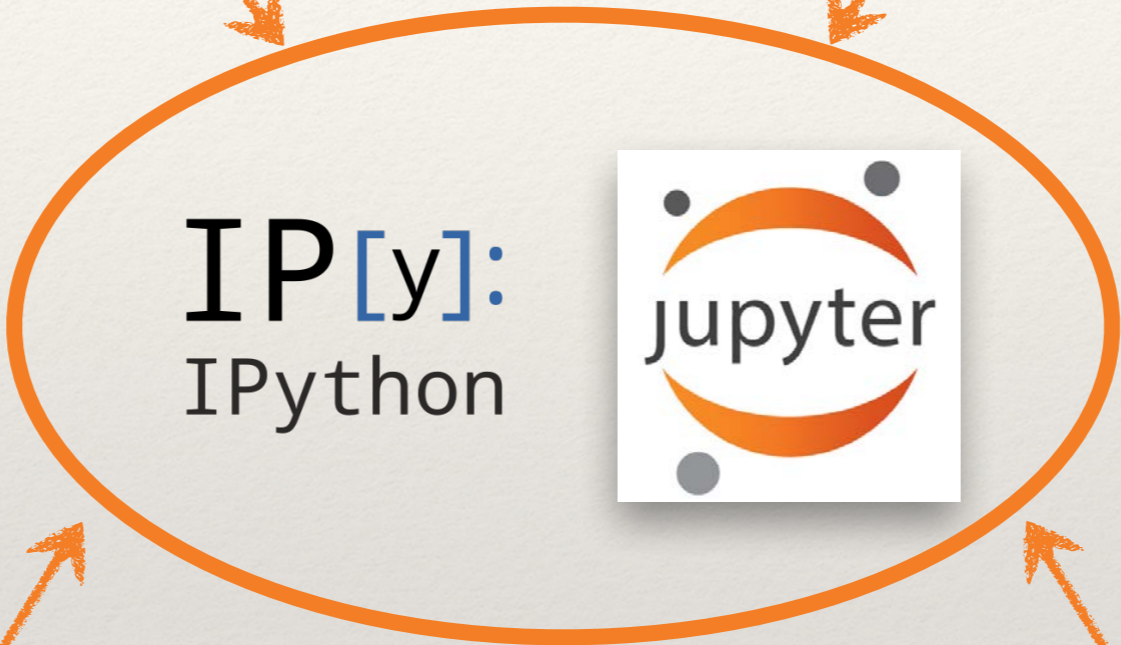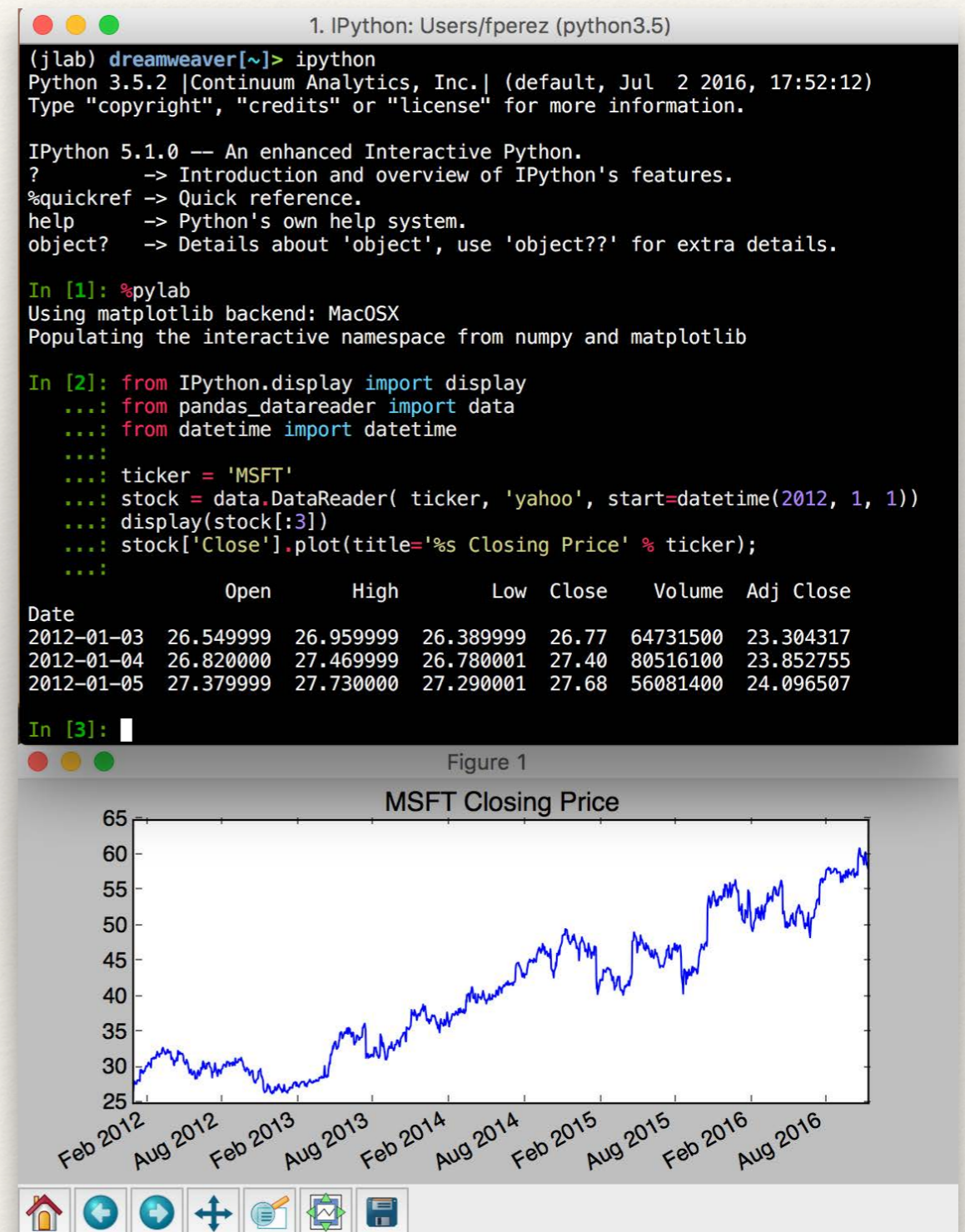Biology: Sequencing
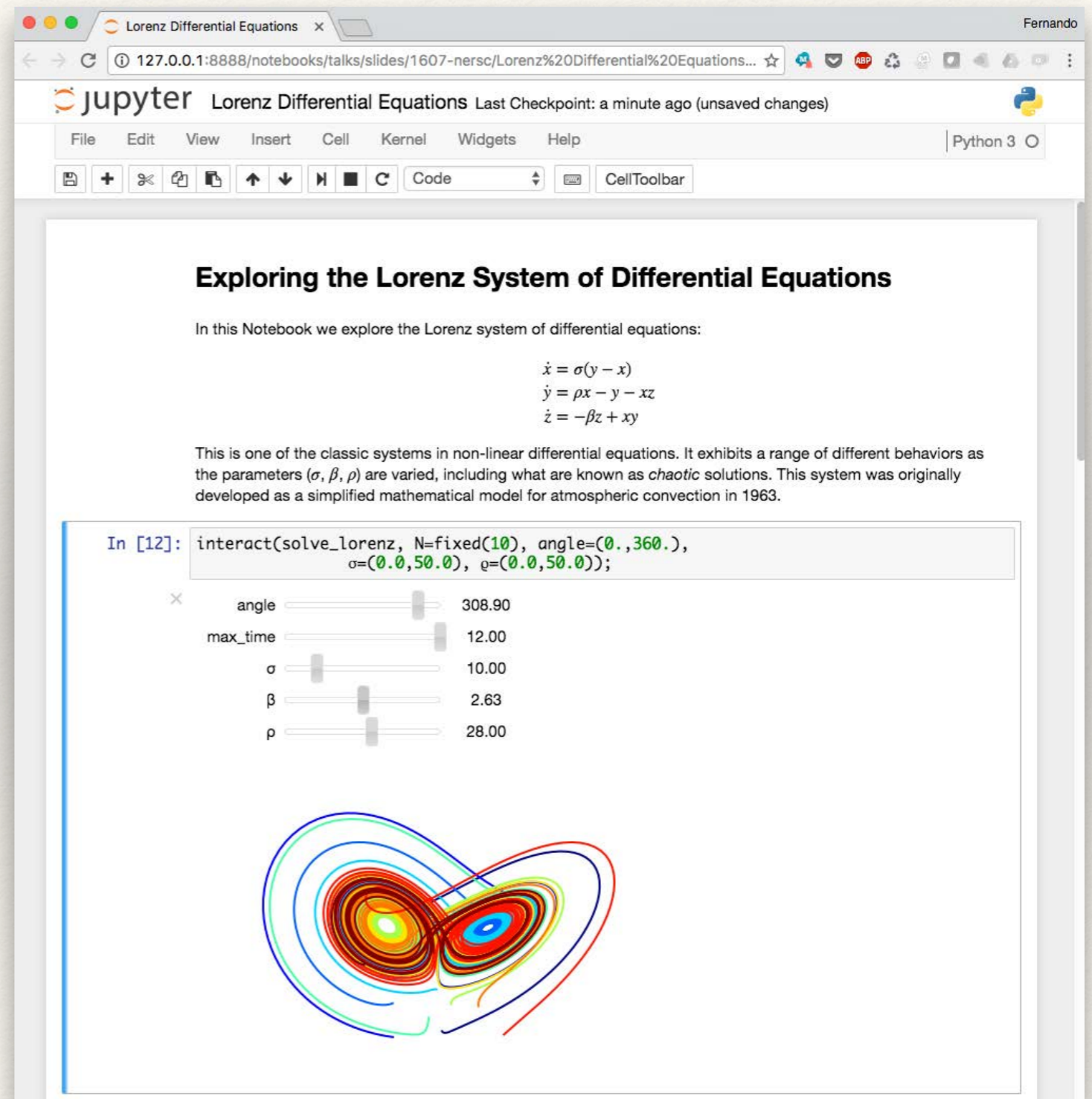


Economics: POS terminals



Neuroscience: EEG, fMRI

IP[y]:
IPython    jupyter

# IPython: Interactive Python, 2001

- Object Introspection (TAB!)

- OS Integration

- Rich terminal client

- GUI support (plots, ...)

- %magic commands

- Embeddable

# The IPython/Jupyter Notebook

- ❖ Rich web client

- ❖ Text & math

- ❖ Code

- ❖ Results

- ❖ Share, reproduce.
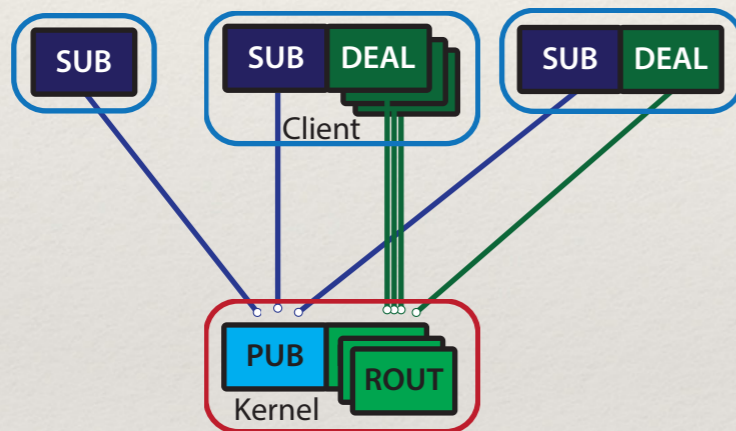
# Funding and partnerships

# Core ideas of the web: HTTP & HTML



HTTP: protocol to connect clients and servers
HyperText Transport Protocol

HTML: format to represent content
HyperText Markup Language

Image credit: eviltester.com

# Core ideas of Jupyter

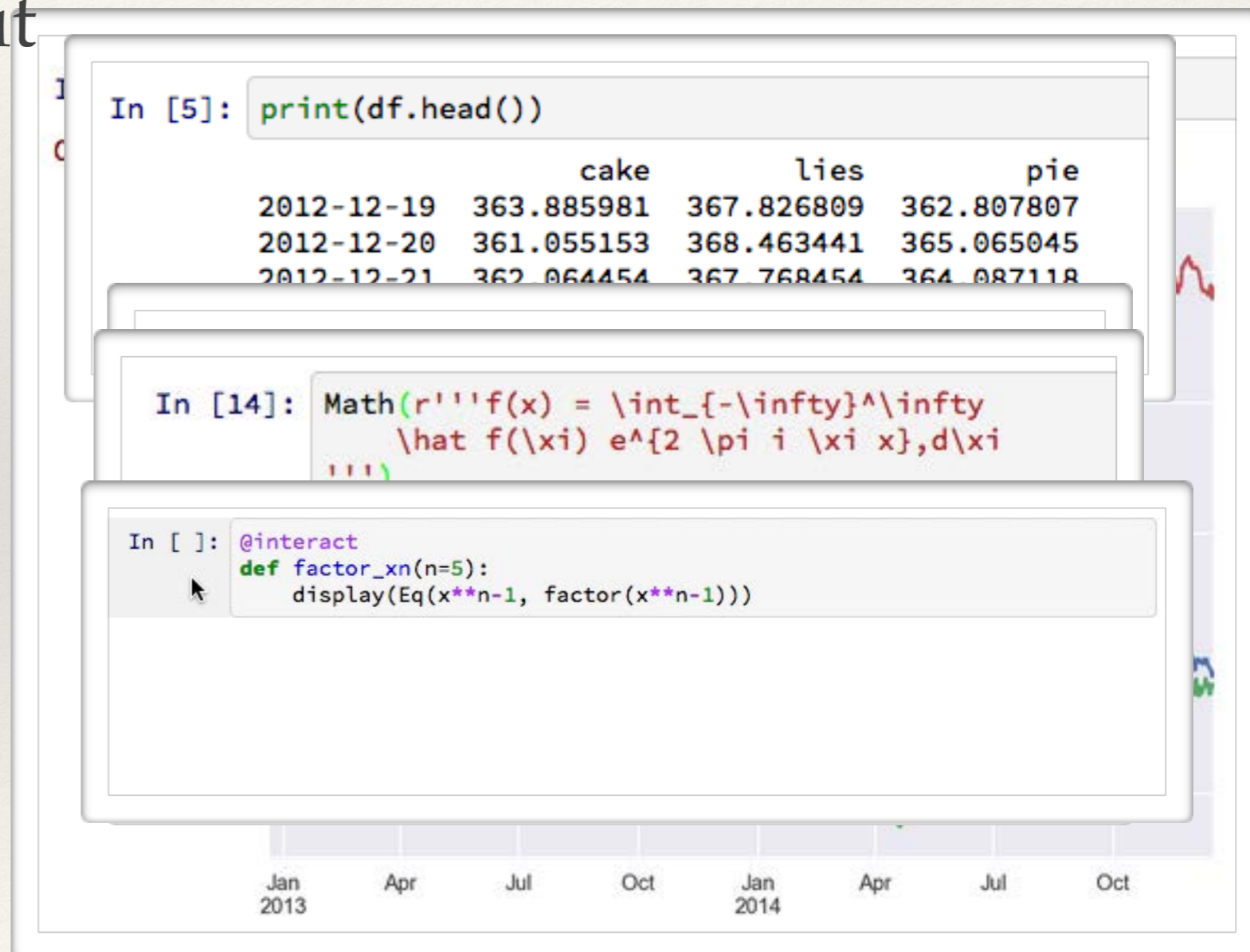## Interactive Computing Protocol



ØMQ + JSON

## Document Format

# Jupyter Protocol
## capture the process of interactive computing

any mime-type output

❖ text

❖ svg, png, jpeg

❖ latex, pdf

❖ html, javascript

❖ interactive widgets

# Jupyter Protocol
# is language agnostic



~75 different kernels: https://github.com/ipython/ipython/wiki/IPython-kernels-for-other-languages

# Notebook: a *data structure*

```
{
  "cells": [
    { ... }, // 3 items
    {
      "cell_type": "markdown",
      "metadata": {},
      "source": [
        "In this Notebook we explore the Lorenz system of differential equations:\n",
        "\n",
        "$$\n",
        "\\begin{aligned}\n",
        "\\dot{x} & = \\sigma(y-x) \\\\\n",
        "\\dot{y} & = \\rho x - y - xz \\\\\n",
        "\\dot{z} & = -\\beta z + xy\n",
        "\\end{aligned}\n",
        "$$\n",
        "\n",
        "This is one of the classic systems in non-linear differential equations. It exhibits a range of different behaviors as",
        "the parameters ($\\sigma$, $\\beta$, $\\rho$) are varied."
      ]
    },
    { ... }, // 3 items
    { ... }, // 3 items
    { ... }, // 5 items
    { ... }, // 5 items
    { ... }, // 5 items
    { ... }, // 3 items
    { ... }, // 5 items
    { ... }, // 3 items
    {
      "cell_type": "code",
      "execution_count": 5,
      "metadata": {},
      "outputs": [
        {
          "data": {
            "image/png": "iVBORw0KGgoAAAANSUhEUgAAAb4AAAAEuCAYAAADx63eqAAAABHNCSVQICAgIfAhkiAAAAAlwSFlz\nAAALEgAACxIB0t1+/AAAIAB...
```

Raw | Parsed

# Reproducible Research

An article about computational science in a scientific publication is **not** the scholarship itself, it is merely **advertising** of the scholarship. The actual scholarship is the complete software development environment and the complete set of instructions which generated the figures.

*Buckheit and Donoho, WaveLab and Reproducible Research, **1995***

# Nature: "the advertising"



Gross, Andrew M., et al. Nature genetics 46.9 (2014): 939-943.

# Notebooks on Github: the "actual scolarship"

# Reproducible Research (2012):
# Paper, Notebooks and Virtual Machine

# [mybinder.org](mybinder.org)



github.com/freeman-lab

github.com/andrewosh

Andrew Osheroff's SciPy'16 talk:
https://www.youtube.com/watch?v=OK6M4w7LYIc

# Gravitational waves detected on Jupyter!



From LIGO Open Science Center, binder-ified: github.com/minrk/ligo-binder

# LIGO: Open Science with Jupyter

# The future of reproducible science?

# Global scientific output doubles every nine years

**Academic Pain**
@AcademicPain

Follow

Trying to keep up with the latest academic literature in your area #highered

# Who is reading the literature?



Figure 2. Percentage of papers needed to obtain 20%, 50% and 80% of the citations received using a two-year citation window, by field, 1900–2005

Larivière & Gingras, arxiv.org/0809.5250

# The scientific literature, today

We are conflating two things:

1. Communication of ideas for others to build upon (hence, reproducibility)

2. Professional credit

*NATURE* | NEWS FEATURE

# Does it take too long to publish research?

**Scientists are becoming increasingly frustrated by the time it takes to publish a paper. Something has to change, they say.**

**Kendall Powell**

10 February 2016

# The literature will be read by the machines



LIGO GW150914 analysis as Jupyter Notebook. 1,000,000+ of these on

Let's "publish" less so we can read more!

# What if…

- All our daily work was captured in a way the machines could read…

- annotated with rich metadata…

- natural language, code, results and data all linked…

- easy for the machines to mine for discovery and credit…

- and less frequent highlights were written in long form, also backed by their "real scholarship" (à la Donoho)?

# What would that look like?

❖ "Executable preprints/blog posts"

  ❖ Capture rapid progress, expose data and software

  ❖ Fully reproducible: build scientific community and knowledge

  ❖ With DOIs - citable as needed.

❖ Peer-reviewed papers:

  ❖ less frequent, high-quality narratives

  ❖ real synthesis of important ideas

**nature** International weekly journal of science

Home | News & Comment | Research | Careers & Jobs | Current Issue | Archive | Audio & V

Archive > Volume 515 > Issue 7528 > Column: World View > Article

NATURE | COLUMN: WORLD VIEW

# Open access is tiring out peer reviewers

As numbers of published articles rise, the scholarly review system must adapt to avoid unmanageable burdens and slipping standards, says **Martijn Arns**.

25 November 2014 | Corrected: 26 November 2014

But in recent months, I received reviews of my own submitted papers that suggest reviewers simply did not read the manuscript properly.

[…]

To protect quality reviewing, a hybrid model should be considered. I suggest a two-tier system, in which some papers are not reviewed before publication at all and are instead subject to a post-publication peer review.

# The "scientific paper of the future"



**Caltech** Library    About    Resources    Services

Envisioning the Scientific Paper of the Future

Monday, January 9, 2017

Location

Caltech; Avery House Dining Hall and Library

Register (requested for catering)
Visitor information

Scientific paper of the Future

Victoria Stodden

Yolanda Gil

Titus Brown

**Living in an Ivory Basement** Stochastic thoughts on science, testing, and programming.

misc    personal    python    science    teaching    testing

The top 10 reasons why blog posts are better than scientific papers

**AGU** American Geophysical Union™

The Geoscience Papers of the Future Initiative

The Geoscience Papers of the Future (GPF) is an initiative to encourage geoscientists to publish papers together with the associated digital products of their research

Data implies software.

Note: This is the second post in a mini-series of blog posts inspired by the workshop Envisioning the Scientific Paper of the Future.

# Some new developments in Jupyter's orbit…

# version control for notebooks?

# nbdime to the rescue!



(**n**ote**b**ook **di**ff and **me**rge: https://github.com/jupyter/nbdime

# JupyterLab: the notebook, evolved…

# The "Notebook"?

# JupyterLab: unifying these ideas



A Collaborative effort:

**Bloomberg**

**CONTINUUM** ANALYTICS

Brian, Jason, Steven, Darian, Sylvain, Carol, Cameron, Farica, Paul, Reese, Kyle, Chris, Ian, Matthias, …

# Live Demo!

Demo credits/thank you:
Brian Granger (Cal Poly SLO)
Jason Grout (Bloomberg)