# Jupyter
## Jupyter in HPC

Feb 28th, 2018

**Matthias Bussonnier**

bussonniermatthias@gmail.com
GitHub: @carreau
Twitter: @mbussonn

# About Me

## Matthias Bussonnier

- A Physicist/Bio-Physicist

- Core developer of IPython/Jupyter since 2012

  - Co-founder, and Steering Council member

- Post doctoral Scholar on Jupyter at BIDS

# Demo(s)

Just to find bugs and make things crash

# Webinar & Outline

- This webinar will be in 3 parts

  - Overview of what is Jupyter + HPC

  - Use case : Suhas Somnath

  - Use case : Shreyas Cholia

- Outline Part 1

  - From IPython to Jupyter

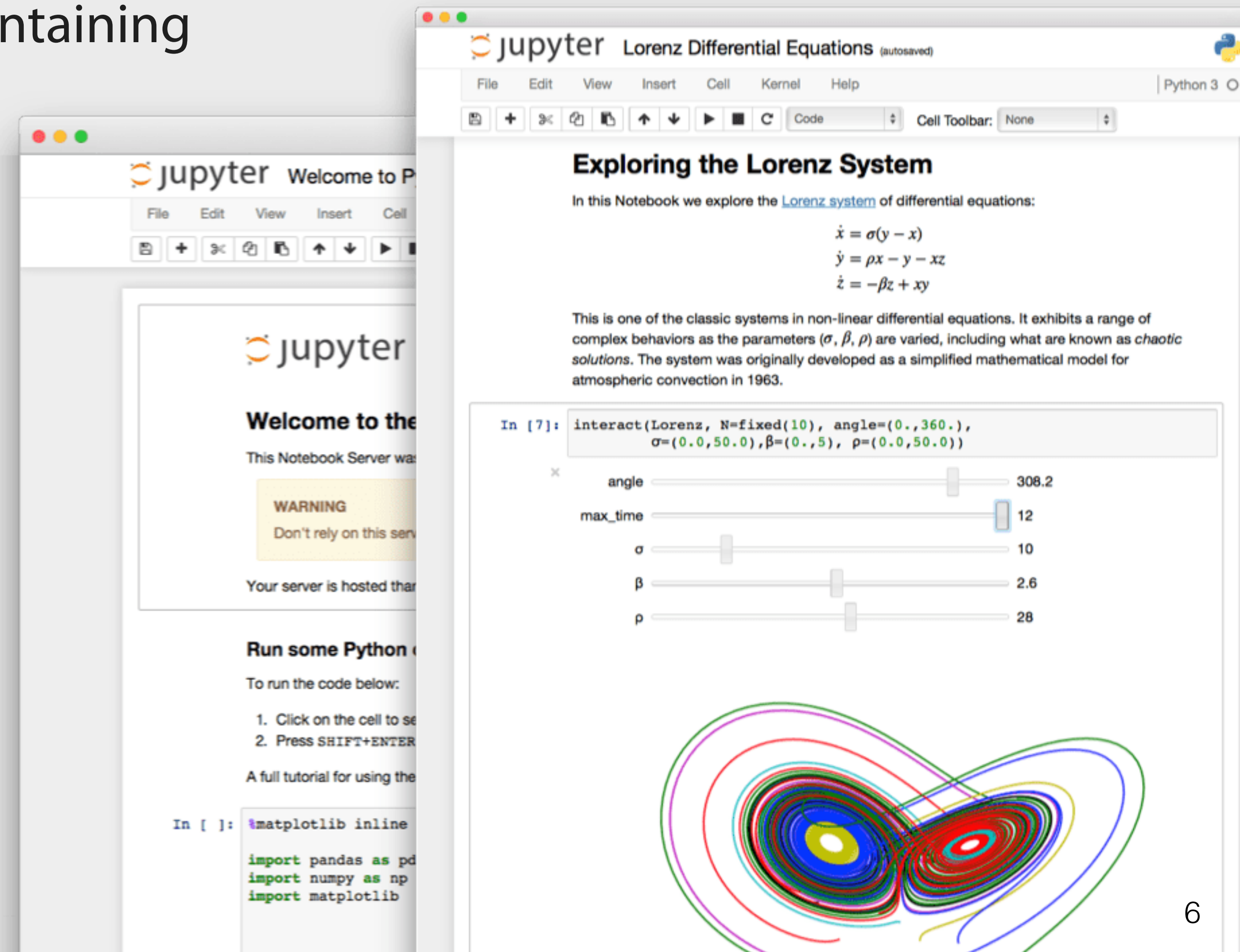  - What is Jupyter

  - Jupyter Popularity

  - Some Jupyter Usage

# From IPython to Jupyter

- 2001: Fernando Perez Wrote "**IPython**"

  - Create IPython for Interactive Python with prompt number, gnu plot integration

  - Replace a bunch on perl/make/C/C++ files with only Python.

- 2011: QtConsole

- 2012: Birth of current **Notebook** (6th prototype)

  - Make IPython "network enabled"

  - Made possible by mature web tech.

- 2013: First non-Python (**Julia)** kernel

- 2014: we **renamed** the Python-Agnostic part to **Jupyter**.

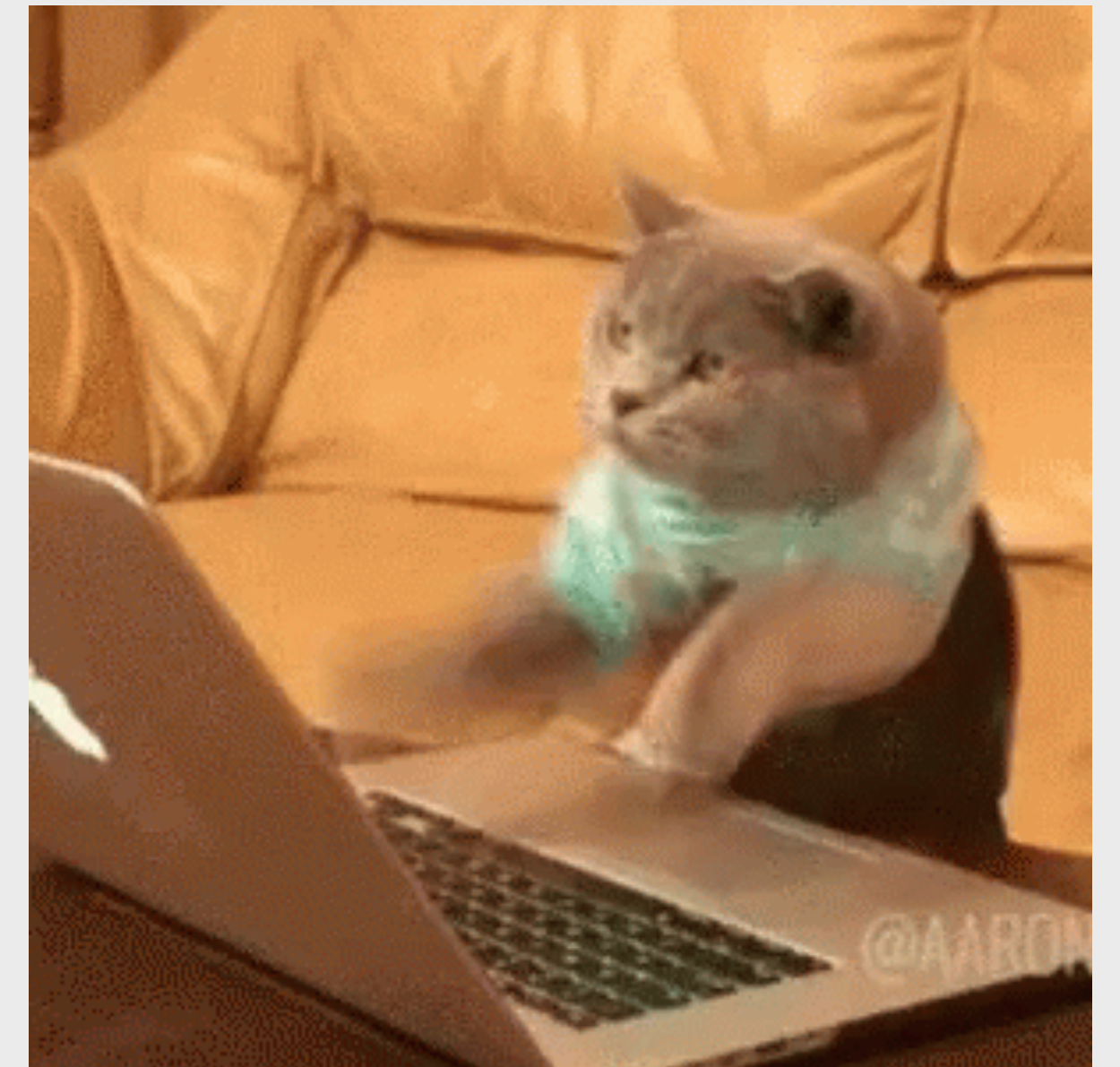- 2018: several millions users & **JupyterLab** released

# What is Jupyter

- Mainly Known for **The Notebook**

  - Web server, a web app, load .ipynb (json), containing

    code, narrative, math and results.

  - Attached to a **Kernel** doing computation.

- Results can be:

  - Static (Image)

  - Interactive (client-side scoll/pan/brush)

  - Dynamic (Call back into Kernel)

# Focused on Exploratory Programming



- IPython was **designed** for exploratory programming, as

  a **REPL** (Read Eval Print Loop) and grew popular, especially

  among scientist who loved it to **explore**.

"IPython have weaponized the tab [completion] key"
– Fernando Pérez

# Open Organisation

- Organisation with Open Governance (https://GitHub.com/jupyter/governance)

- Funded by Grants and Donations, and Collaborations

# Protocols and Formats

- Jupyter is also a set of **Protocols and Formats** that reduce the **N-frontends $\times$ M-backends** problem to a **M-Frontends + N-backends**,

  - Open, Free and Simple.

    - JSON (almost) everywhere

    - Notebook document format,

    - Wire protocol

  - Thought for Science and **Interactive** use case.

    - Results embedded in documents no "Copy past" mistake.

    - Scale from Education to HPC jobs.

# Ecosystem

**Frontends**: Notebook, JupyterLab, CLI, *Vim, Emacs, Visual Studio Code, Atom, Nteract, Juno...*

**Kernels**: Python, *Julia, R, Haskell, Perl, Fortran, Ruby, Javascript, C/ C++, Go, Scala, Elixir... 60+*

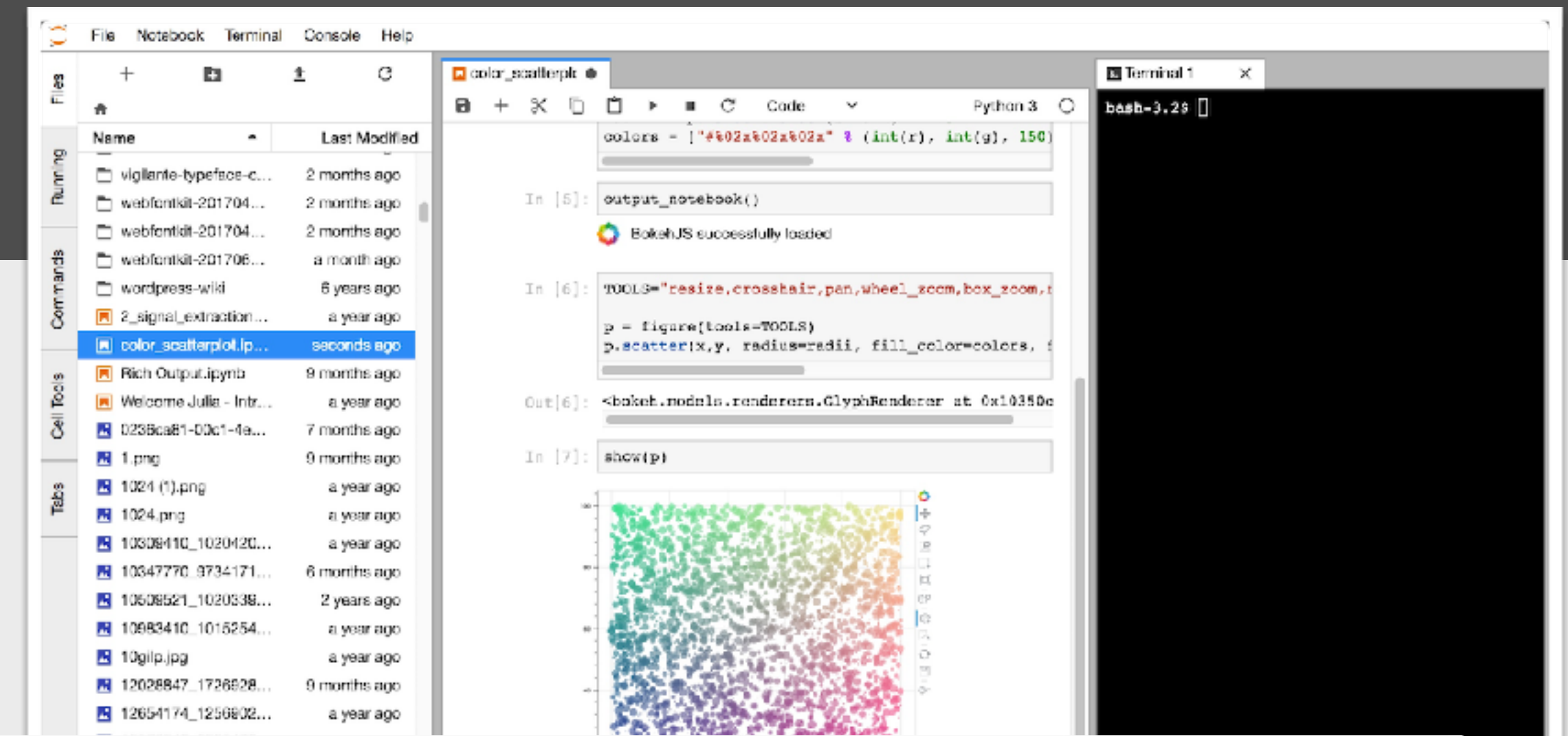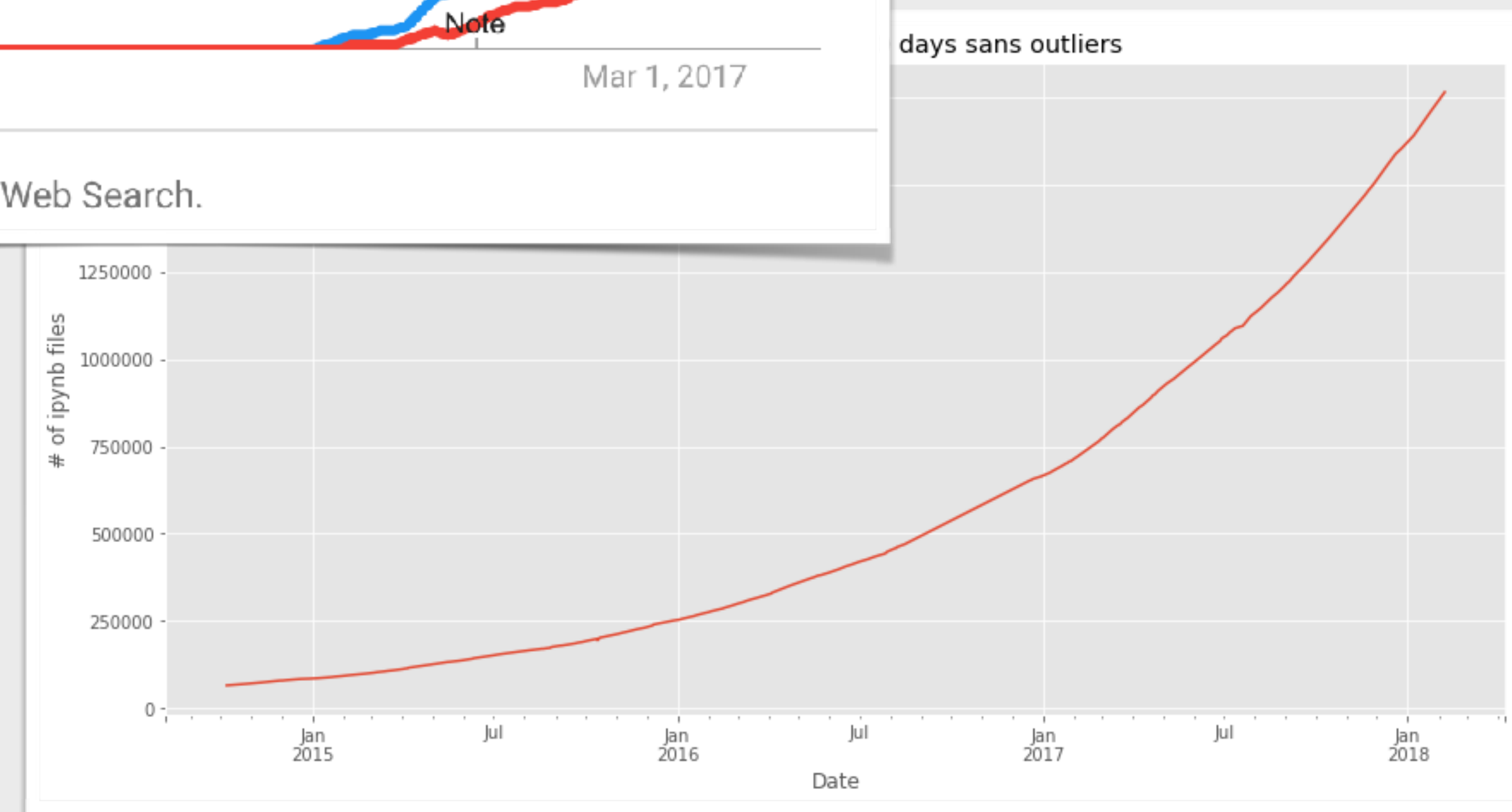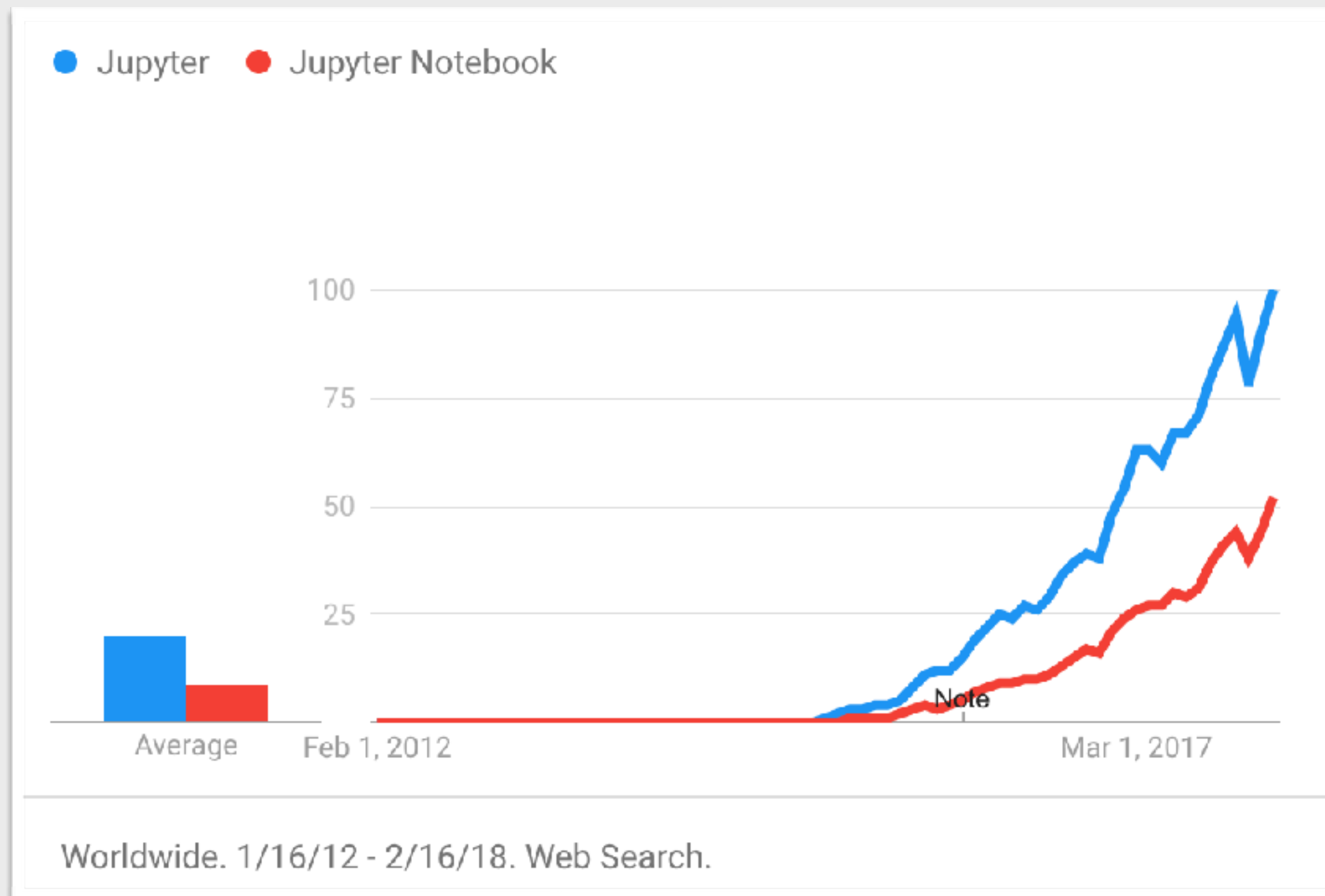**Building Blocks:** Nbformat, JupyterHub, Kernel Gateway...

# JupyterLab



- Extends the notebook interface

  with text editor, shell, ...etc

- is it and IDE ?

- If by I you mean Interactive,

  then yes

# Popularity



Worldwide. 1/16/12 - 2/16/18. Web Search.



https://github.com/parente/nbestimate

# Interactivity

**Popularity**

- Coding is not the end goal of most of our users. A simple, single tool, with friendly interface helps.

- Persisting kernel state allows to iterate only on part of an analysis.

- Notebook interface give the interactivity of the REPL with the edit-ability and linearity of a script with intermediate result. Aka "Literate Computing"

# Separation of states

- Computation, narrative/visualisation in different processes.

    - Robust to crashes

    - Can "Share" and analysis / notebook without having to "rerun"

    - Trustworthy (No copy-past issues).

- Cons:

    - Understanding that document/kernel can have different states can be challenging.

    - Notebook format is not as widespread as others.

Popularity

# Popularity

# Network enabled / web based

- User love fancy colors and things moving. Using D3 and other

> **Bojan Marković**
> Feb 20
>
> You'll only take Spyder from my cold, dead… Oooooh, pretty shiny colors, inline graphics.. Does it come in fuchsia? :)

  - dynamic libraries are highly popular

  - Usable by novices and power-users

  - Users w/ different expertise (Numerical Methods, Visualization,...)

- Seamless transition to HPC: Kernel Menu > Restart on Cluster

- Document persist if code crash.

- Can be Zero-Installation (See JupyterHub).

  - A web browser is all you need.

# JupyterHub

- Multi-users Jupyter deployment

  - Not (Yet) Realtime collaboration

- Each user can get their own process/version(s)/

  configuration(s)

  - Hooks into any Auth

  - Only requires a browser

- Not limited to running Jupyter (e.g. work with RStudio,

  OpenRefine...)

# Use Cases

## HPC

- Batch Jobs

  - You can run notebook "headless"

  - Parametrized notebook as "reports" you can interact with later

- Interactive Cluster.

  - Run a Hub (hook into LDAP/PAM...)

  - Run notebook servers on a Head node

  - Run Kernels on head Node/fast queue

  - Extra Workers (e.g. dask) on Batch queue/cluster.

# HPC deployment



Head Node

Head Node

Hub

Auth

Server

Kernels

Head Node

Cluster:
- Dask
- MPI
- IPyparallel

# HPC Misconceptions

**We need to run JupyterHub on the cluster:** No

- Hub, Server, Kernels, (and Workers):
  - Do not have to be on the same machine
  - Do not have to use the same environment

**A Kernel is a (single) language:** No

- A Kernel is a preconfigured computation environment. It can be:
  - A queue, a hardware resource (GPU, SSD...), A location (like a beam-line)
- Example of Python, Cython, Julia, R, Fortran, Rust, C calling each other in same notebook

**Every User have the same environment:** No/No

- Kernels and notebook server can be configured independently
- Subset of users could use different server versions w/ different extensions.

**JupyterHub is Limited to Jupyter:** No

JupyterHub Can run RStudio, Open Refine.

# Danger !

Despite Notebook being great, some limitations:

Most if not all **document state** is **in your browser** !

- Watch out for flaky network connections !

- Do Not close your Laptop Lid*/Tab* !

Workaround:

- Wrap computations (especially long), in Futures

- Use Caching.

Interrupting in compiled code is hard.

Large outputs/notebooks can crash the browser

# Some Jupyter Usage



**Ligo**



**Pangeo**



**Cern's SWAN**

# Ligo

- Some events analysis with Jupyter

- Subset of data + env put online

- Run the analysis yourself on Binder[1] and listen to the waves





[1] https://github.com/minrk/ligo-binder

# Pangeo (pangeo-data.github.io)

- Effort from Atmosphere / Ocean / Land / Climate (AOC) science

  community

- unified effort

- Cloud based

- Recent Technologies

  - Dask, Jupyter



Matt Rocklin Blog post on pangeo-data.github.io

# Cern Swan (swan.web.cern.ch)

- Share platformed for Data Analysis

- Sync W/ $HOME directory

- 0-install

- Share Data

- Provide example gallery with 1-click-fork

# The Shape of Things to come

Classic Notebook -> JupyterLab transition

- Stabilisation

- Transfer of extensions

- Collaboration:

    - Google retired Real-time API

    - Who "executes" problem

- Long Running Jobs

JupyterHub

- Horizontal (and Vertical) Scaling

- Audits APIs  (Hippa Compliance ?)

- "Federation" (binder) / Intercommunication

# CFP- Ends March 6th

# Question(s)
# while we change
# speakers ?