# First Notebook War

## PyData DC 2018

Martin Skarzynski

2018-11-17

# About me

- Cancer Prevention Fellow
  - at National Cancer Institute
- Co-Chair, Bioinformatics & Data Science
  - at Foundation for Advanced Education in the Sciences
- Website: https://marskar.github.io
- Twitter: @marskar

# Title inspired by a tweet by Philip Guo

# Previous conflicts

- Spaces versus Tabs
- Emacs versus Vim
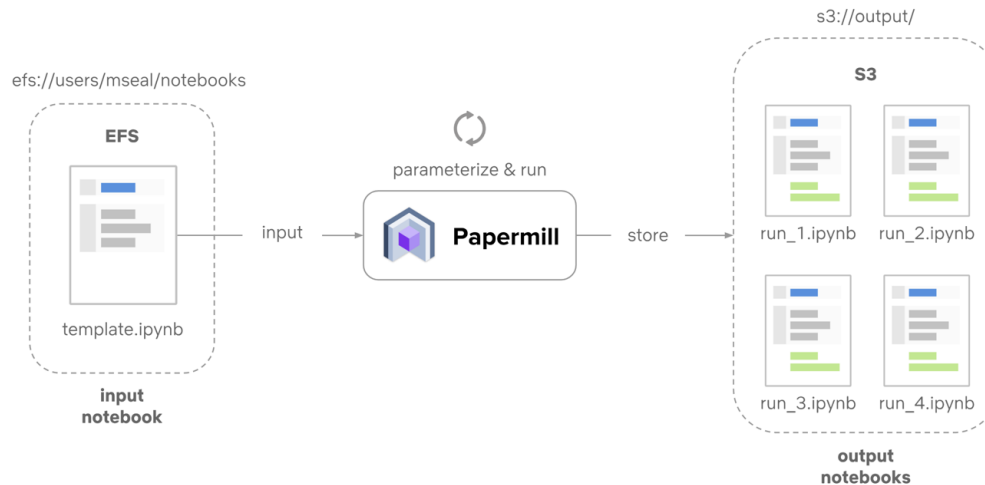


- Source: https://goo.gl/images/U2KpcG

# Spaces versus Tabs

- Code editor setup: Tab = 4 spaces
- GNU Make requires tabs!
- Use spaces, get paid more!
    - according to blog post by David Robinson (@drob)'s

**Salary differences between developers who use tabs and spaces**
From 12,426 professional developers in the 2017 Developer Survey results, who provided tabs/spaces and salary

# Notebooks

Jupyter notebooks

- are data science tools
- built on IPython by Fernando Perez (@fperez)
- combine Markdown text, code, and output
- help data scientists communicate goals, methods, and results
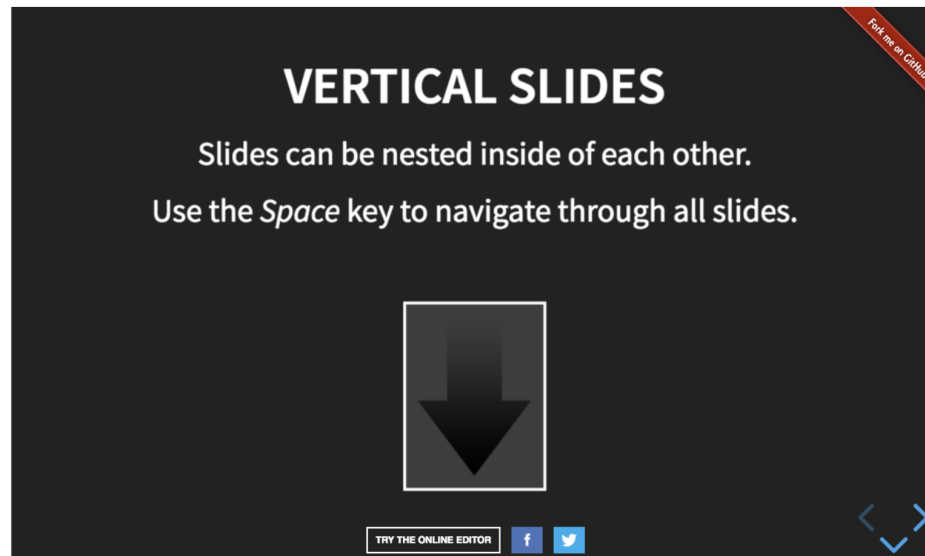- used in academia, Amazon, Netflix, and PayPal

# Joel doesn't like notebooks

- "I Don't Like Notebooks" by Joel Grus (@joelgrus) at JupyterCon 2018
- Slides
- Video
- Modularity and Reusability

# Joel also doesn't like vertical slides

**PLEASE COME TO MY NEXT TALK:**
**"I DON'T LIKE THAT SLIDESHOW PROGRAM WHERE SOMETIMES THE NEXT**
**SLIDE IS RIGHT AND SOMETIMES THE NEXT SLIDE IS DOWN"**



@joelgrus #jupytercon

# Tim Hopper (@tdhopper) likes notebooks

# DataFramed

- Podcast by Hugo Bowne-Anderson (@hugobowne)
- Episode 44 with Brian Granger (@ellisonbg)
- JupyterLab



« We are entering an era where large, complex organizations need to scale interactive computing with data to their entire organization in a manner that is collaborative, secure, and human centered. »
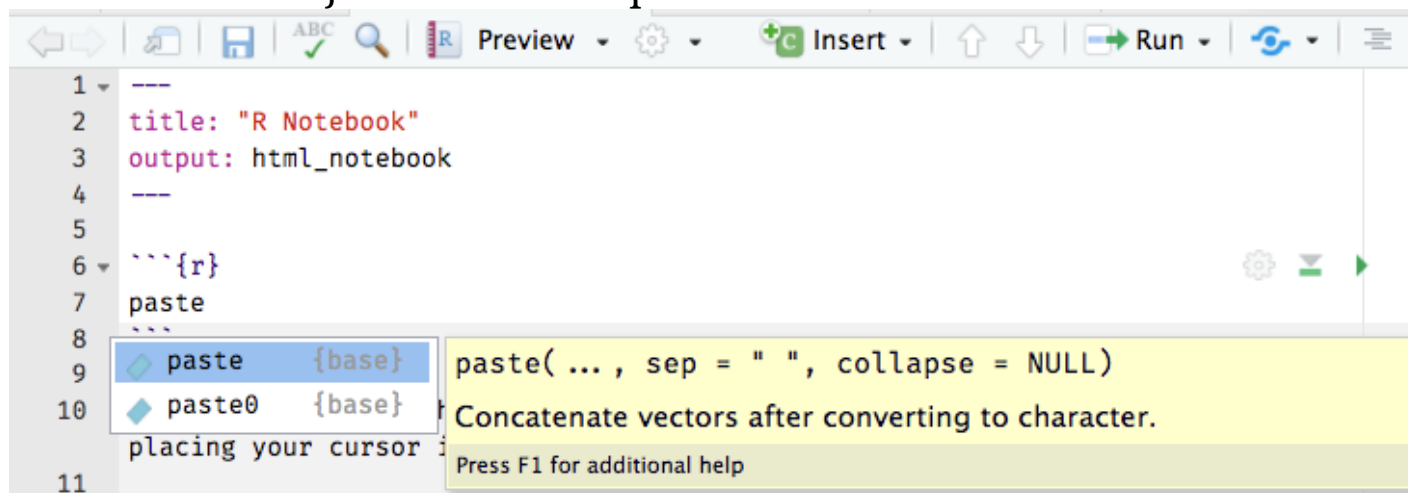
BRIAN GRANGER

DataFramed BY DataCamp

# Yihui Xie (@xieyihui)'s Blog post

- I used R markdown to make these slides!
- No problems with version control
- R notebooks are just another output format!

# Why use notebooks

- Literate programming
- Rendered by GitHub and nbviewer
- Google colab
- Kaggle kernels
- Binder

# Notebook tools

1. version control tool for notebooks - nbdime
2. work with Jupyter notebooks and scripts in parallel using JupyText
3. configure notebooks to run on markdown (md) files with notedown
4. create and run Jupyter notebooks from scripts and md files with nbless



Traditional Jupyter notebook

# Write modules!

- Imports
    1. Standard Library
    2. Third Party
    3. User Defined
- Definitions
    - Classes
    - Functions (for more check out Steven Lott's PyData DC tutorial)
- Type Hints
- Docstrings (with examples!)
- `if __name__ == '__main__':`
- Function call(s), e.g. `doctest.testmod(verbose=True)`
- doctest: docstring examples -> test suite (with `unittest` API)
- run test suite with `unittest` or `pytest`
- use cookiecutter for project structure
- deploy projects/packages to PyPI

# Thanks for listening!