

Always Evolving:

A wild foray into the world of genomics and visualization

By Blake Joyce & Asher

And there'll be Plotly coding later on

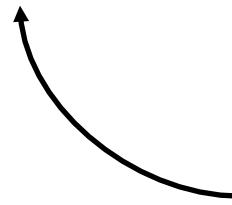
Let's Talk Biology

I promise it's not as boring as your Intro Bio class probably made it

Besides you'll be human all your life (hopefully): it's useful to know how you work

► We propose:

Genetics has always been a big data & visualization problem



Did he just say “job security”?

The 3 Challenges of Genetics: in brief

Need to visualization three major things in particular:

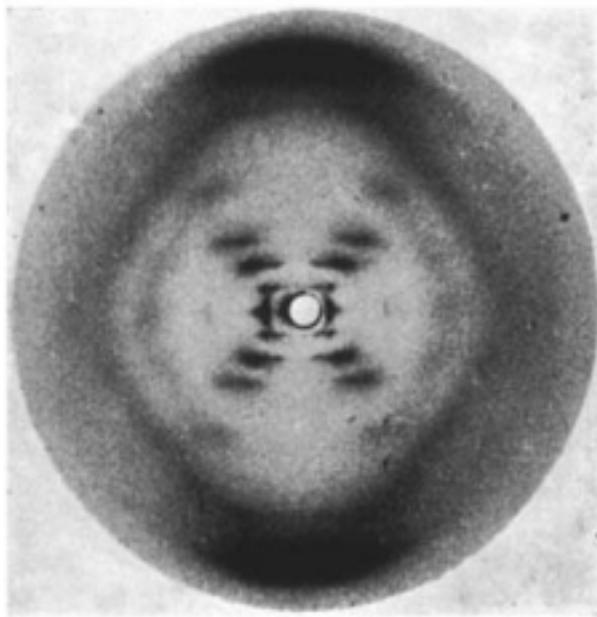
1) Initial visualization of genetic sequence (a.k.a. ‘Sequencing’)

- 1.1) How to “see” chemical sequences at A, T, G, C
- 1.2) Crazy sizes of data (hard to get all that information)

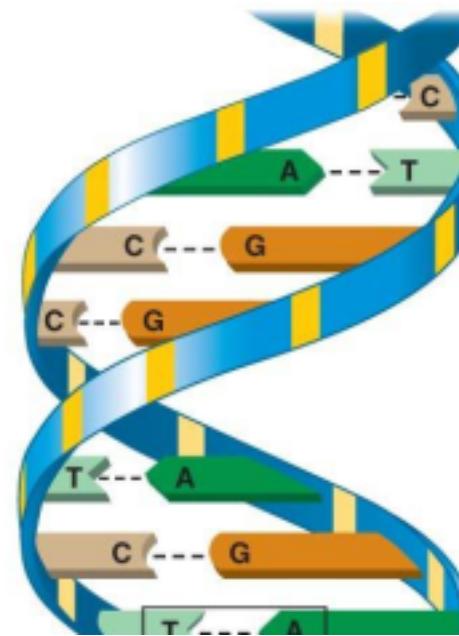
2) Visualization to understand the sequence

- 2.1) Hard to ‘comprehend’ that amount of data
 - 2.2) Compare to other genetics to understand evolution, gene function, etc
- 3) Assess and visualize to teach/conceptualize genetics (another conversation)

Challenge 1.1: How Can We See the Sequence?

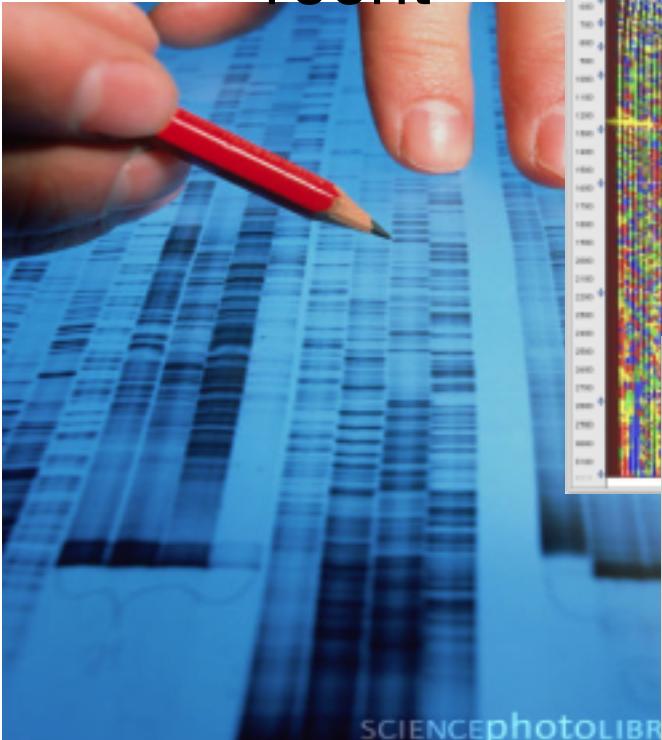


=

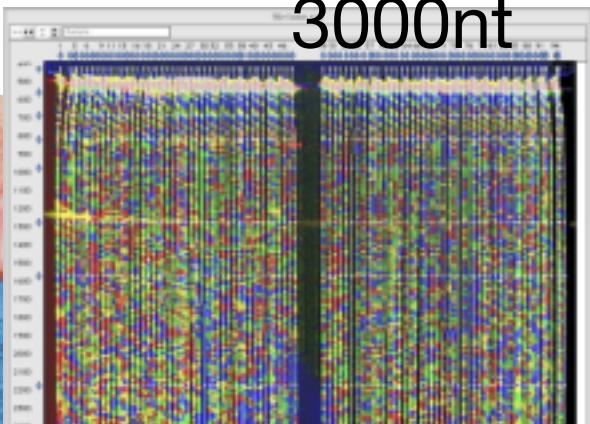


DNA sequencing

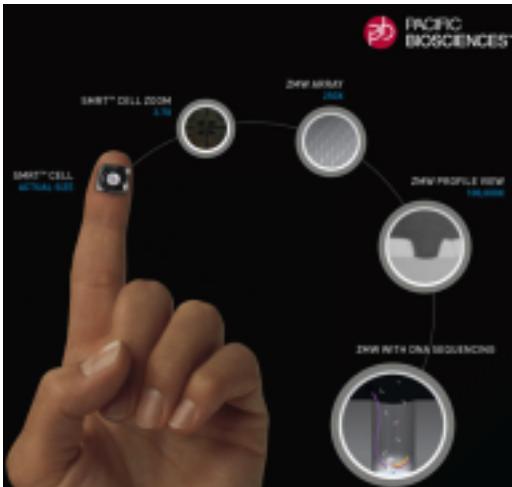
100nt



3000nt



10-20 kb reads



512,000,000nt/2 weeks



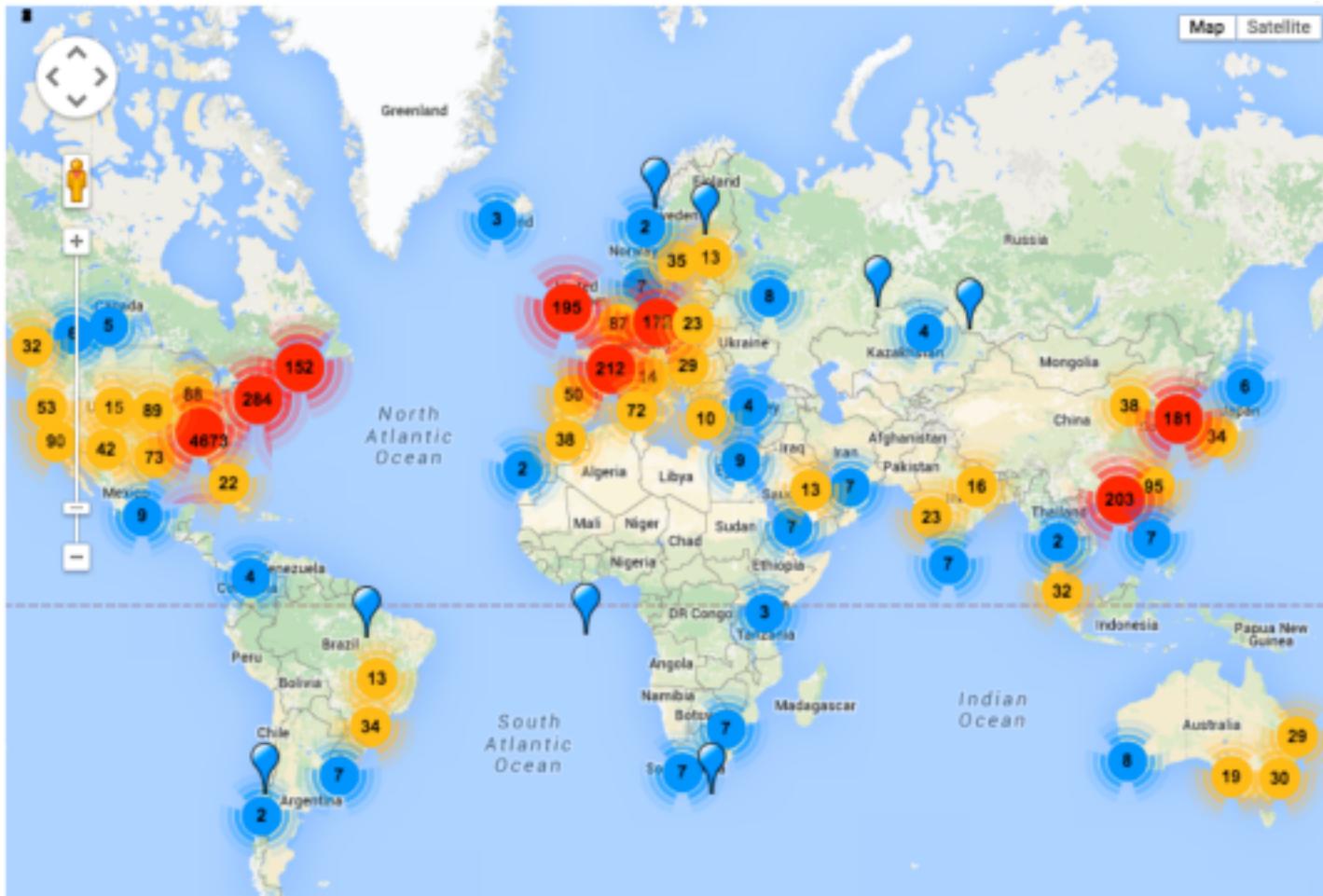
Illumina HiSeq X Ten



In < 3 Days:

- Generates 1800 GB of raw data
- 1,600,000,000,000 bases of DNA sequence
(500 human genomes)
- +18,000 humans/year @ ~\$1000/each

And there are sequencers across the globe



Cool. So We Can See It Pretty Well.

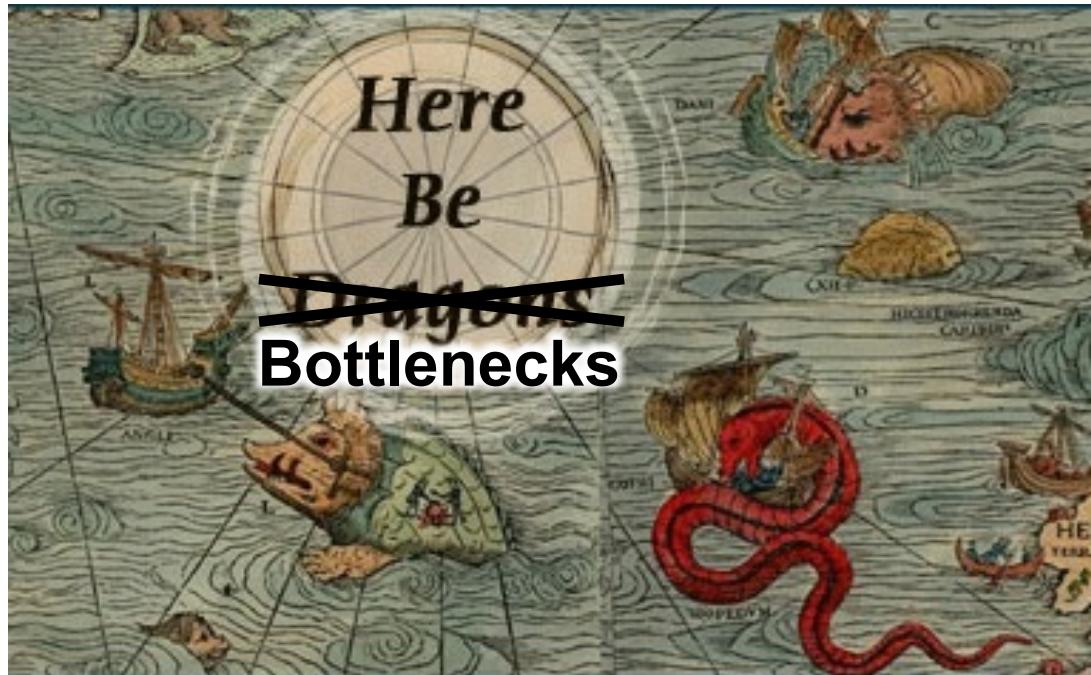
Challenge 1.1 solved?

The short answer is kind of....

Long answer: Challenge 1.1 isn't a bottleneck any more

Challenge 1.2: Size Does Matter ←(scientifically proven)

- “Here be bottlenecks”
- Data science and cyberinfrastructure to the rescue!
- Got to move, analyze, assess, and keep forever



And there's a lot of data. (feel free to say: "How much is there?" in unison)

Genome size as HPE (Harry Potter Equivalents)



Plant Genome Sizes

Arabidopsis: 37 HPES

Rice: 100 HPES

Maize: 650 HPES

Wheat: 6000 HPES

Paris japonica: 50,000 HPES

One half of one human genome.

HPES

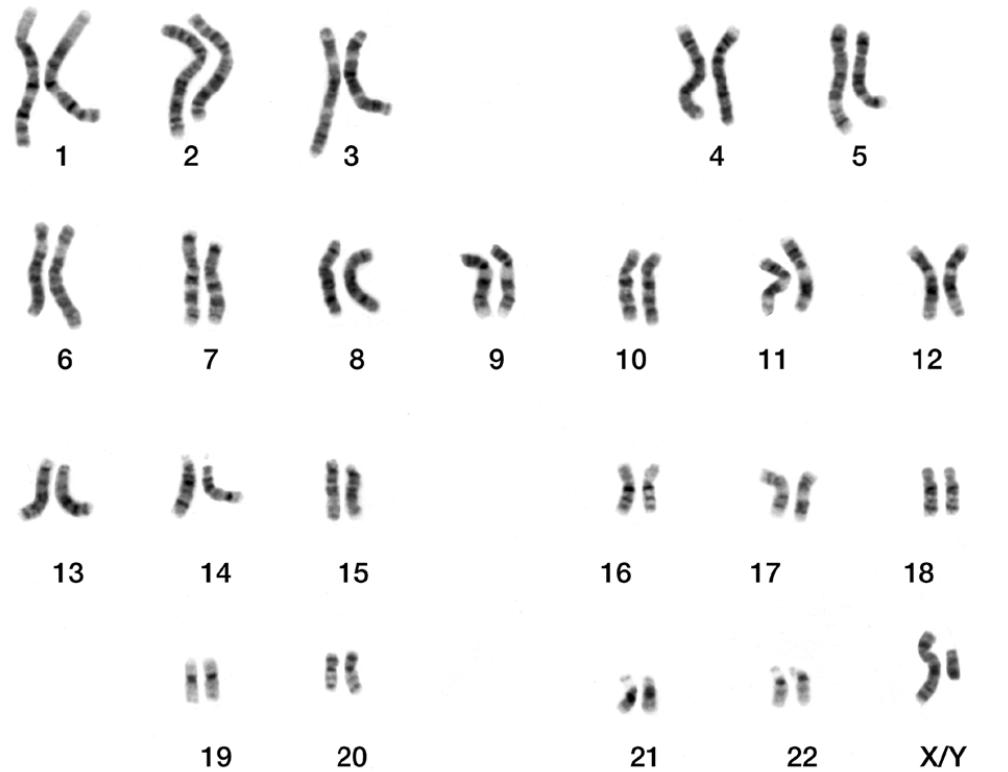


Homo sapien

And that's just half...

You get two
halves...

2128 HPE (>6 billion ATGC's)



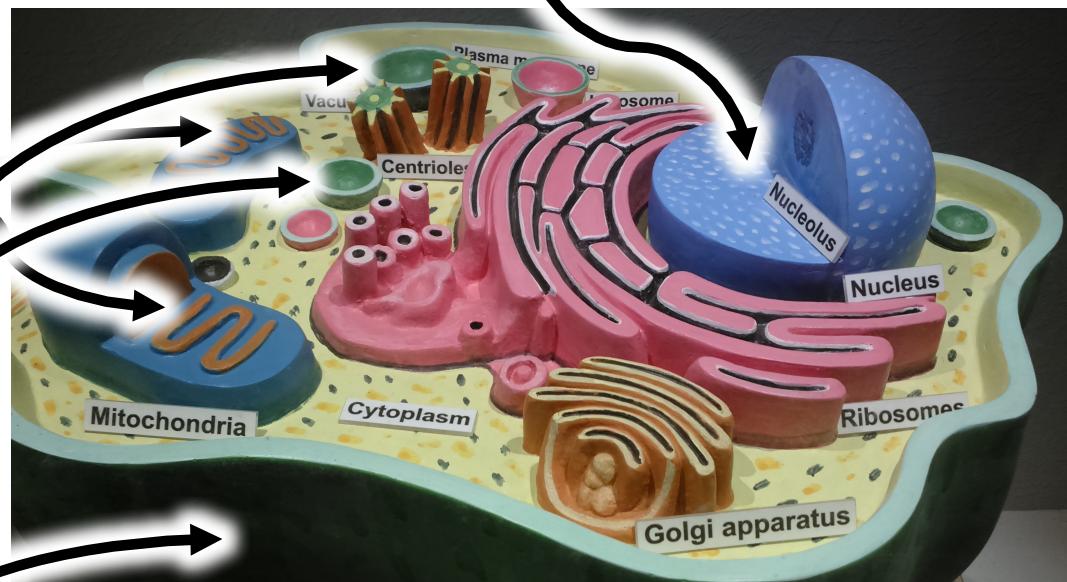
And that's just the nucleus...

There's DNA in other places in the cell...

But DNA Also Be Here

And chloroplasts (plants)

Here Be the “Two Halves”



Speaking of visualization problems: this is NOT how a cell looks

And that's just one cell...

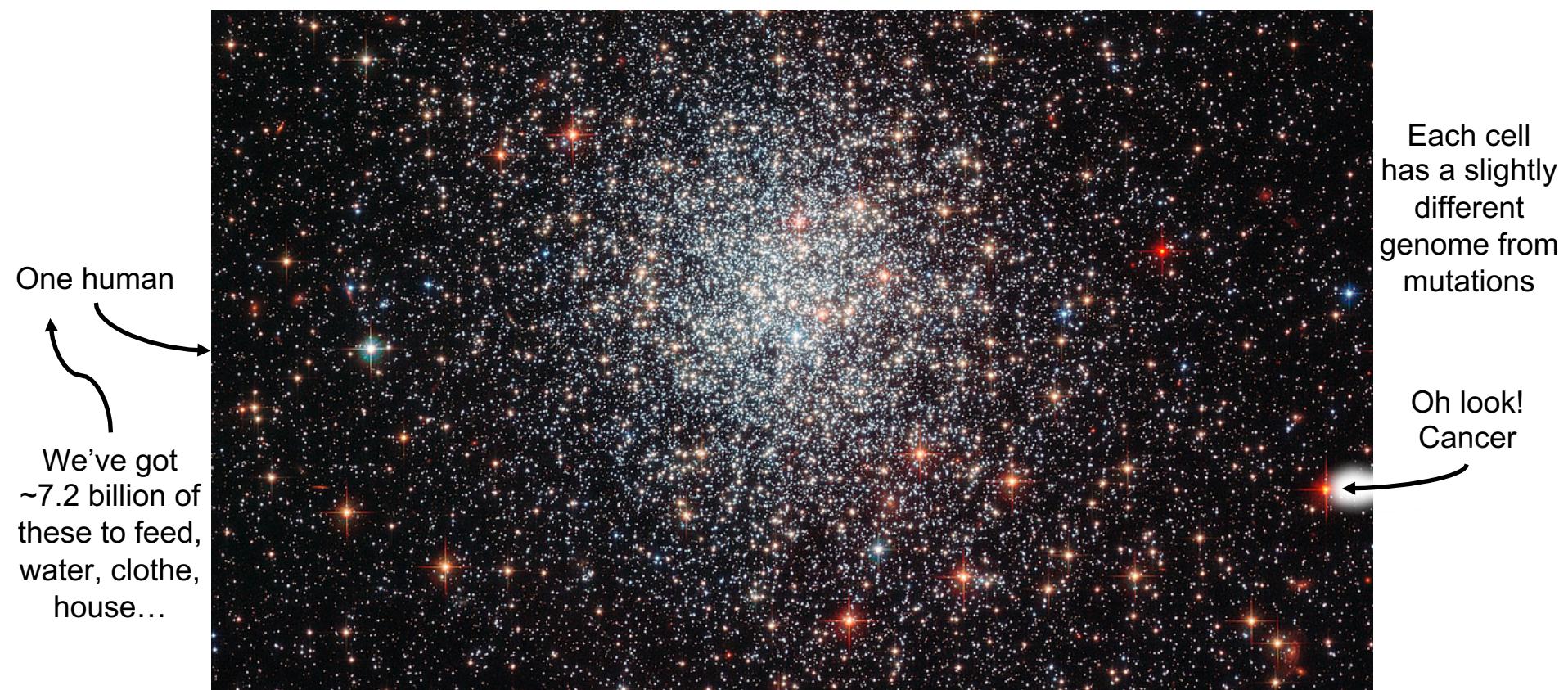
You've got trillions of cells working together in a cacophony of concerted harmony to maintain your life (biology) every moment of every day...

I'm not sure how to visualize this



Trillions of Your Cells + Trillions of Microbes

Here's the only way I can think to visualize that:



And that's just genomes...

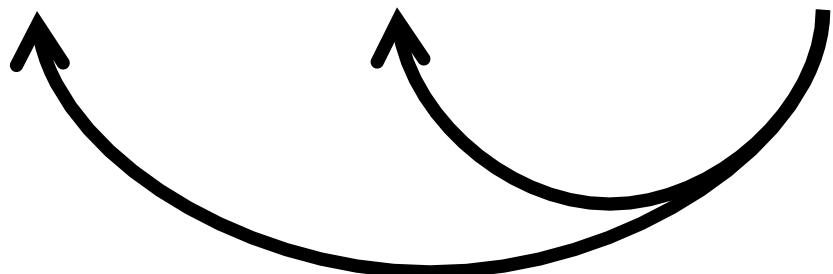


“Textbook Biology”

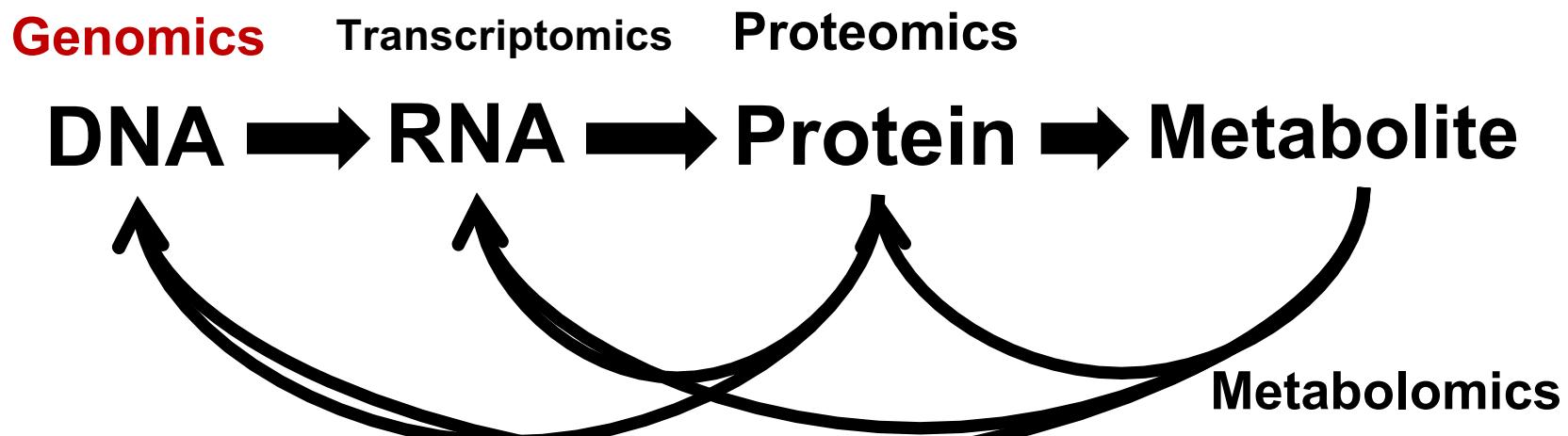
Challenge 1.2 Mark II: integrate all the **genome** data with other data types

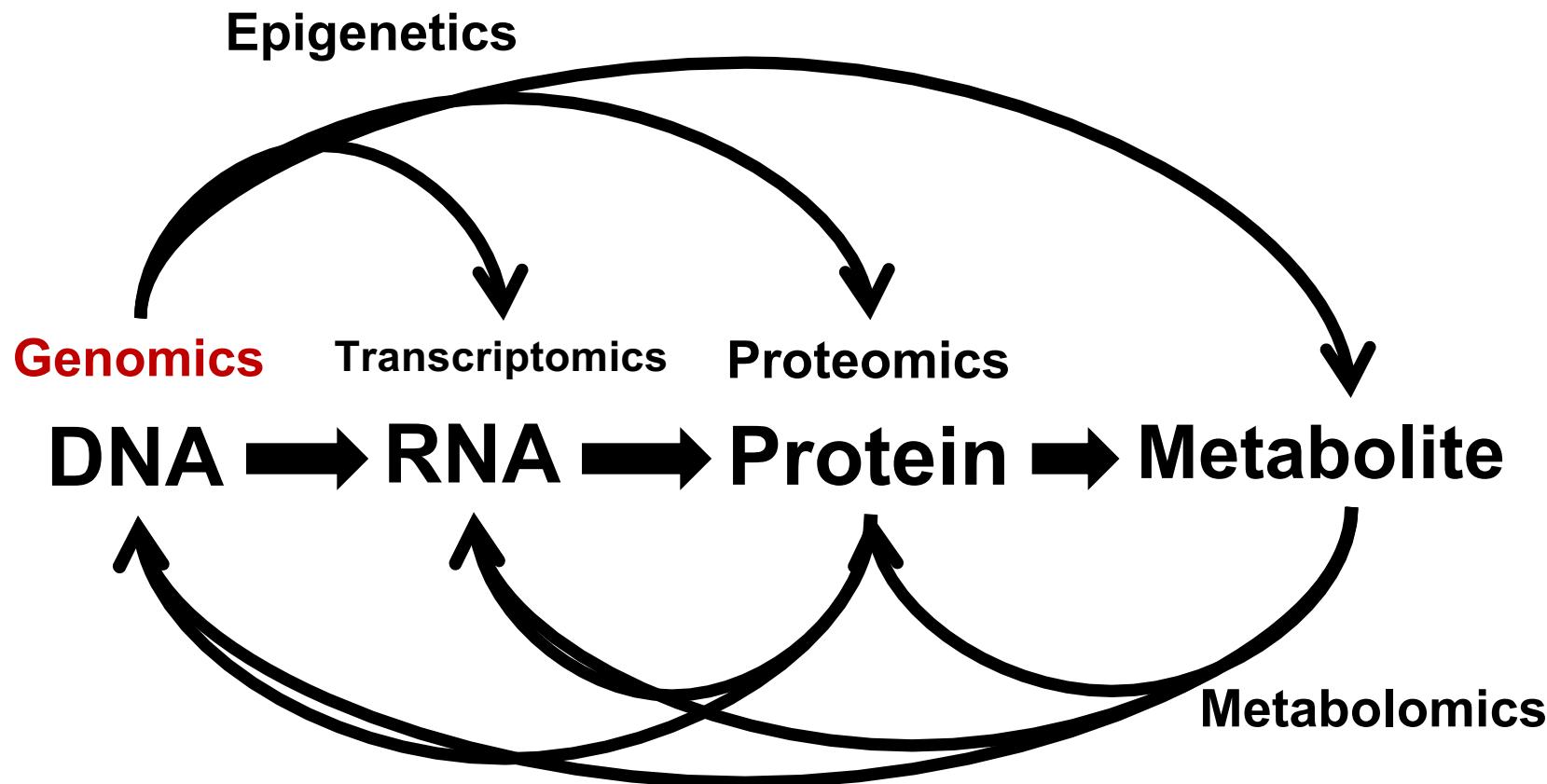
Genomics Transcriptomics Proteomics

DNA → RNA → Protein → Metabolite

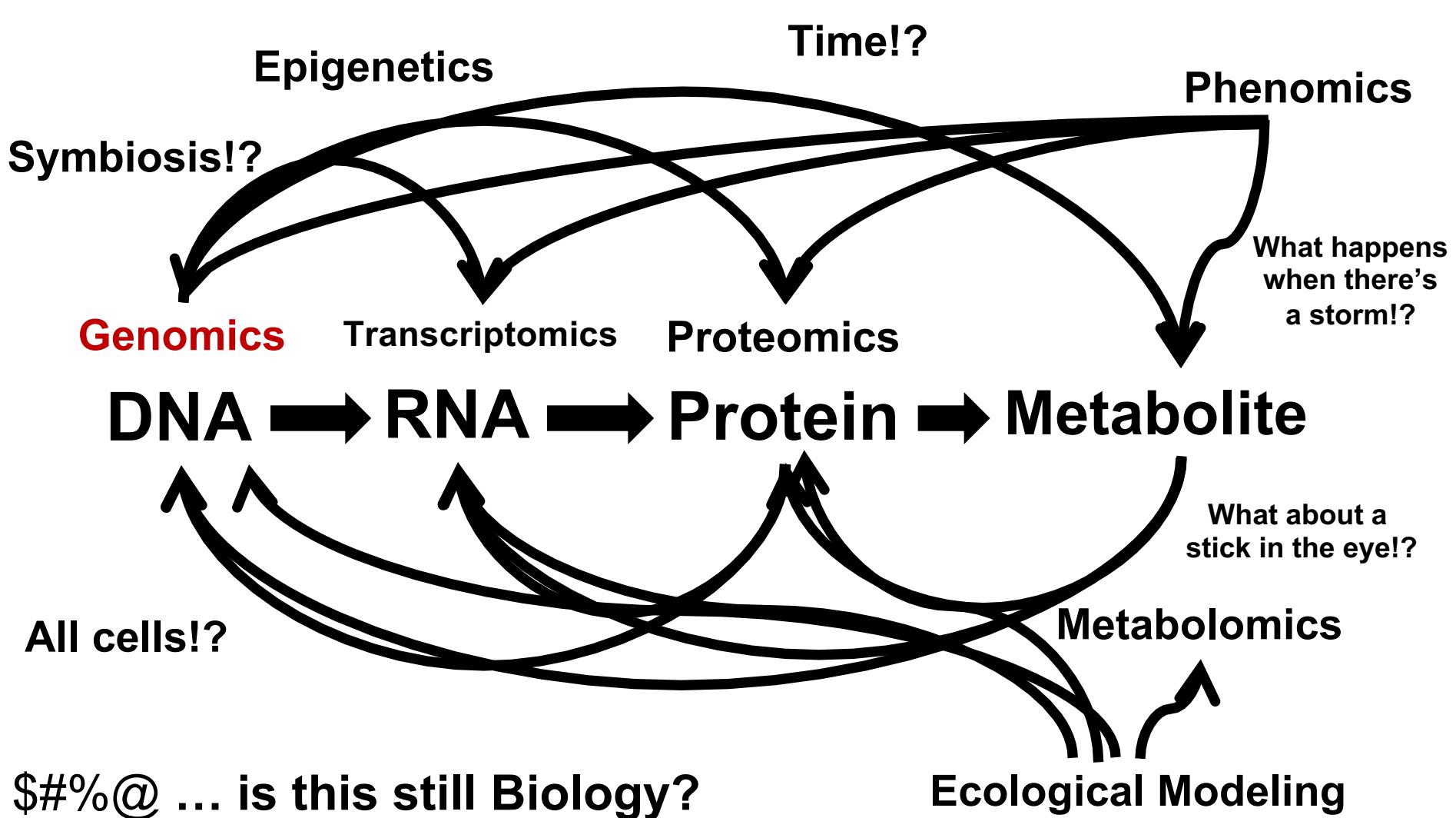


Truer “Textbook Biology”





Even Truer Biology. We're still working on the textbooks.



DON'T PANIC



Embrace the chaos!
The effort is worth it...

Challenge 1.2 Solved: Sign up for CYVERSE



I am a CyVerse SI Team member: think of me as your chaos sherpa (Beaker is my spirit animal)

Shameless Plug:



PhTea

PHARMACY | SCIENCE | AGRICULTURE | ENGINEERING | MEDICINE

Uniting Worlds



Solving Problems

Weekly on Tuesdays 8-10 AM BIO5 Institute

HACKY HOUR

PHARMACY | SCIENCE | AGRICULTURE | ENGINEERING | MEDICINE



DRINK TOGETHER. WORK TOGETHER.

Weekly on Thursdays 4-7 PM Gentle Ben's

TACCCCTTCAGGAACAAATAGGATGGATGACAAATAATCCACCTATCCCAGTAGGAGAA
ATTTATAAAAAGATGGATAATCCTGGGATTAAATAAAATAGTAAGAATGTATAGCCCTA
CCAGCATTCTGGACATAAGACAAGGACCAAAGGAACCCTTAGAGACTATGTAGACCG
GT...
GA
GA
CC.
AT
GC.
AT
GG

So We've Got Data...

AG
GA
CC
AT
GC
AT
GG

And We Can Move It and Fun Things Like That

AGCCAACAGCCCCACCAGAACGAGAGCTTCAGGTCTGGGGTAGAGACAAACACTCCCC
TCAGAACGAGGAGCCGATAGACAAAGGAACGTATCCTTAACCTCCCTCAGGTCACTC
TTTGGCAACGACCCCTCGTCACAATAAGATAGGGGGGCAACTAAAGGAAGCTCTAT
TAGATACAGGAGCAGATGATACAGTATTAGAAGAAATGAGTTGCCAGGAAGATGGAA

GA
GA
CC.
AT
GC.
AT
GG

But what does the genome look like? (to us)

AGCCAACAGCCCCACCAGAAGAGAGCTTCAGGTCTGGGGTAGAGACAACAACTCCCCC
TCAGAACGAGGAGCCGATAGACAAGGAACTGTATCCTTAACTCCCTCAGGTCACTC
TTGGCAACGACCCCTCGTCACAATAAGATAGGGGGGCAACTAAAGGAAGCTCTAT
TAGATACAGGAGCAGATGATACTAGTATTAGAAGAAATGAGTTGCCAGGAAGATGGAA

TACCCCTTCAGGAACAAATAGGATGGATGACAAATAATCCACCTATCCCAGTAGGAGAA
ATTTATAAAAAGATGGATAATCCTGGGATTAAATAAAATAGTAAGAATGTATAGCCCTA
CCAGCATTCTGGACATAAGACAAGGACCAAAGGAACCCCTTAGAGACTATGTAGACCG
GTTCTATAAAACTCTAACAGAGCCGAGCAAGCTTCACAGGAGGTAAAAAATTGGATGACA
GAAACCTTGTGGTCCAAAATGCGAACCCAGATTGTAAGACTATTTAAAAGCATTGG
GACCAGCGGCTACACTAGAAGAAATGATGACAGCATGTCAGGGAGTAGGAGGACCCGG
CCATAAGGCAAGAGT ATTCAAGCTACCATA
ATGATGCAGAGAGGC .GTGTTCAATTGTG
GCAAAGAAGGGCACA AAGGGCTGTTGGAA
ATGTGGAAAGGAAGGACACCAAATGAAAGATTGTACTGAGAGACAGGCTAATTTTTA
GGGAAGATCTGGCCTTCCTACAAAGGGAAAGGCCAGGGATTTCTTCAGAGCAGACAG
AGCCAACAGCCCCACCAGAAGAGAGCTTCAGGTCTGGGTAGAGACAAACAACCTCCCC
TCAGAACGCAGGAGCCGATAGACAAGGAACTGTATCCTTAACCCCTCAGGTCACTC
TTTGGCAACGACCCCTCGTCACAATAAGATAGGGGGCAACTAAAGGAAGCTCTAT
TAGATACAGGAGCAGATGATACAGTATTAGAAGAAATGAGTTGCCAGGAAGATGGAA

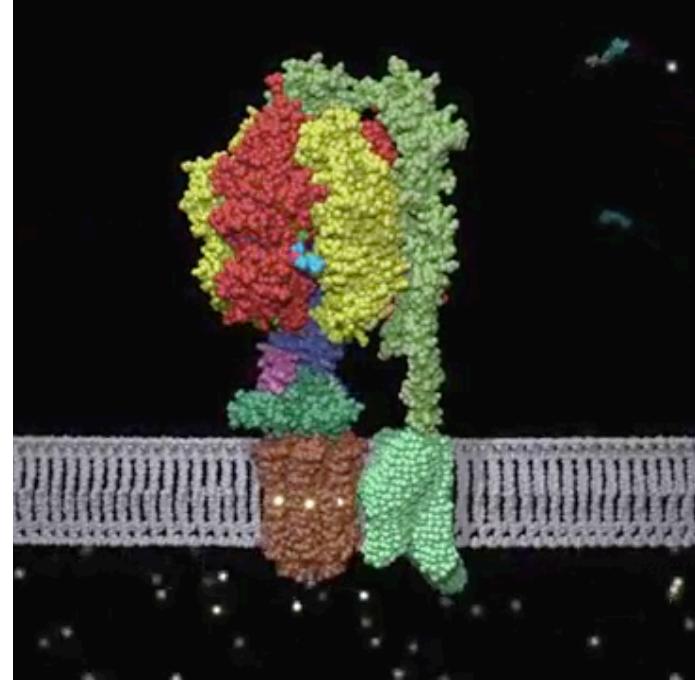
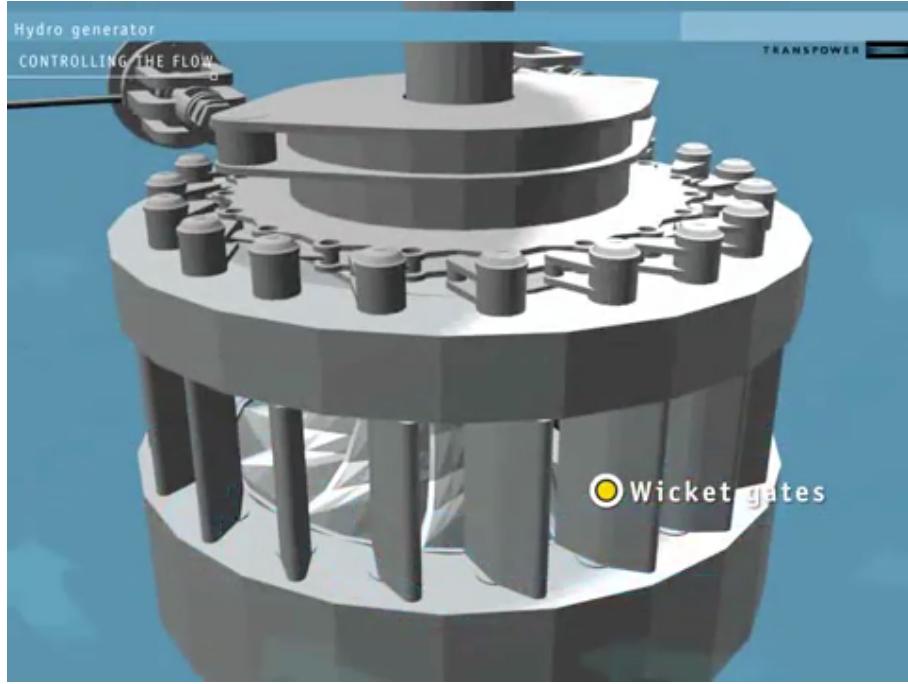
A lot of this.

TACCCCTTCAGGAACAAATAGGATGGATGACAAATAATCCACCTATCCCAGTAGGAGAA
ATTTATAAAAAGATGGATAATCCTGGGATTAAATAAAATAGTAAGAATGTATAGCCCTA
CCAGCATTCTGGACATAAGACAAGGACCAAAGGAACCCTTAGAGACTATGTAGACCG
GTTCTATAAAACTCTAACAGAGCCGAGCAAGCTTCACAGGAGGTAAAAAATTGGATGACA
GAAACCTTGTGGTCCAAAATGCGAACCCAGATTGTAAGACTATTTAAAAGCATTGG
GACCAAC
CCATAAA
ATGATC
GCAAAC
ATGTGC
GGGAAGATCTGGCCTTCCTACAAAGGGAAAGGCCAGGGATTTCAGAGCAGACCAAG
AGCCAACAGCCCCACCAGAAGAGAGCTTCAGGTCTGGGGTAGAGACAACAACTCCCCC
TCAGAACGAGGAGCCGATAGACAAGGAACGTATCCTTAACTCCCTCAGGTCACTC
TTTGGCAACGACCCCTCGTCACAATAAGATAGGGGGGCAACTAAAGGAAGCTCTAT
TAGATACAGGAGCAGATGATACAGTATTAGAAGAAATGAGTTGCCAGGAAGATGGAA

But this turns into cool
micromachines (life)

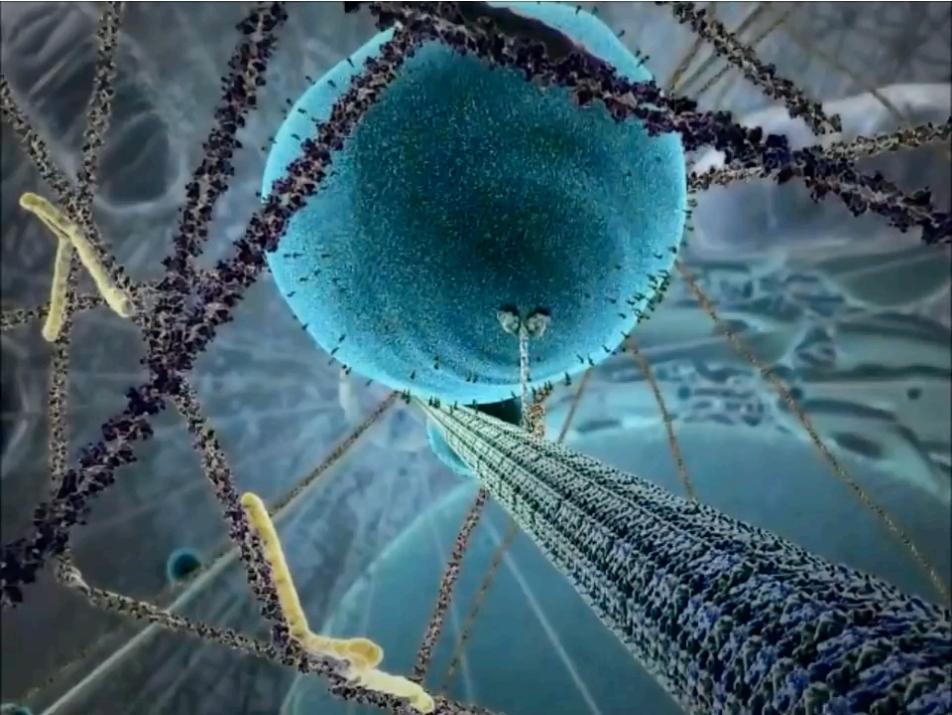
CCCGG
CCATA
TTGTG
GGGAA
TTTTA

And Bio-micromachines (Proteins) Do All Kinds of Things
(and biology did it first)



\ Uses rotation and gradients to convert kinetic energy into other forms

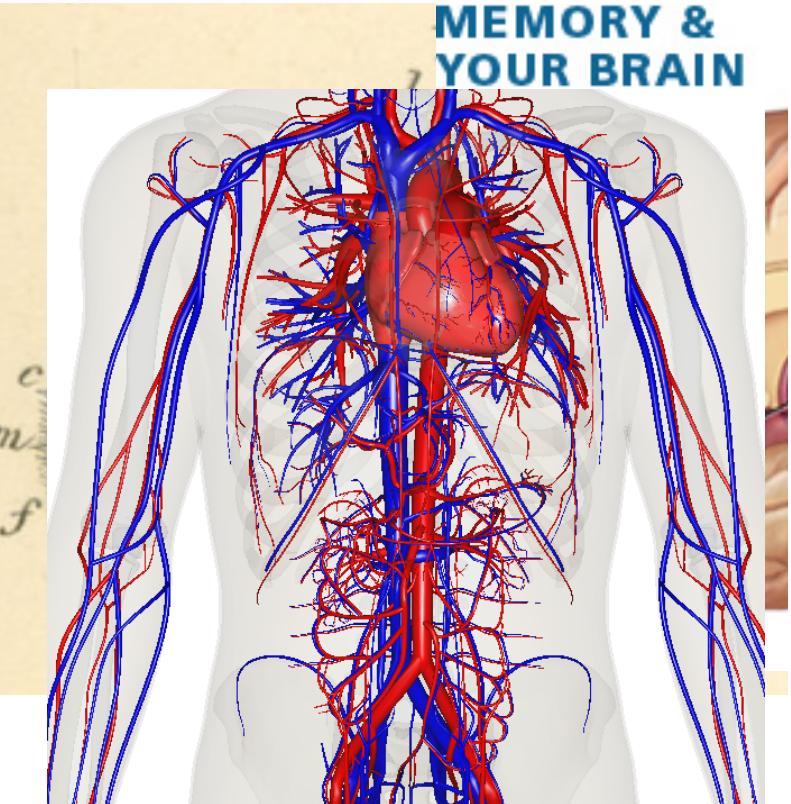
And Bio-micromachines (Proteins) Do All Kinds of Things (and biology did it first)



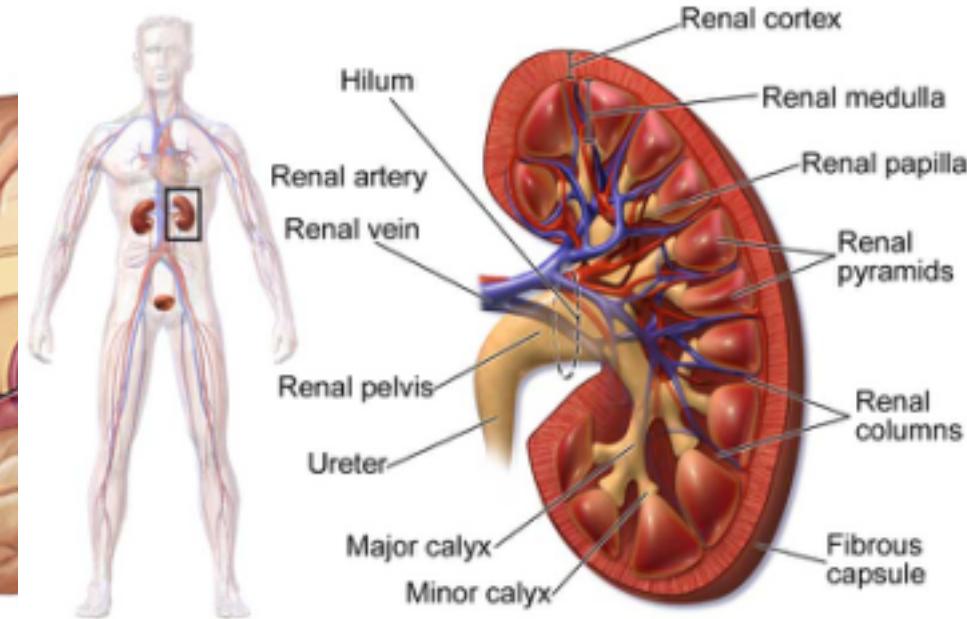
Converts chemical energy into kinetic energy (movement)

And That's Just One Cell

Combinations of cells yield crazy things (that we mostly take for granted):



**MEMORY &
YOUR BRAIN**



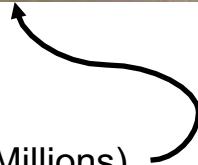
Kidney Anatomy

And That's Just One Cell

Combinations of cells yield crazy things (that we mostly take for granted):



DARPA Project (\$\$\$Millions)



2016 P&G Gymnastics Championships

NBC NEWS.COM
#NIGHTLINE

And That's Just Humans (who frankly are boring)



NAT GEO
WILD
NATGEOWILD.COM

Zombified snails

And There's Some Interesting Applications (like bulletproof skin)



Spider goats



Human skin cells with spider silk

And it's all right here
(somewhere)

GCATTGG
GACCCGG
TACCATA
AATTGTG
GTTGGAA

Challenge 2: How do we understand all this stuff?

Challenge 2: How do we understand all this stuff?

CC. Need to transform the raw data
GT

Need to transform the raw data

We've Come a Long Way

Hand comparisons!

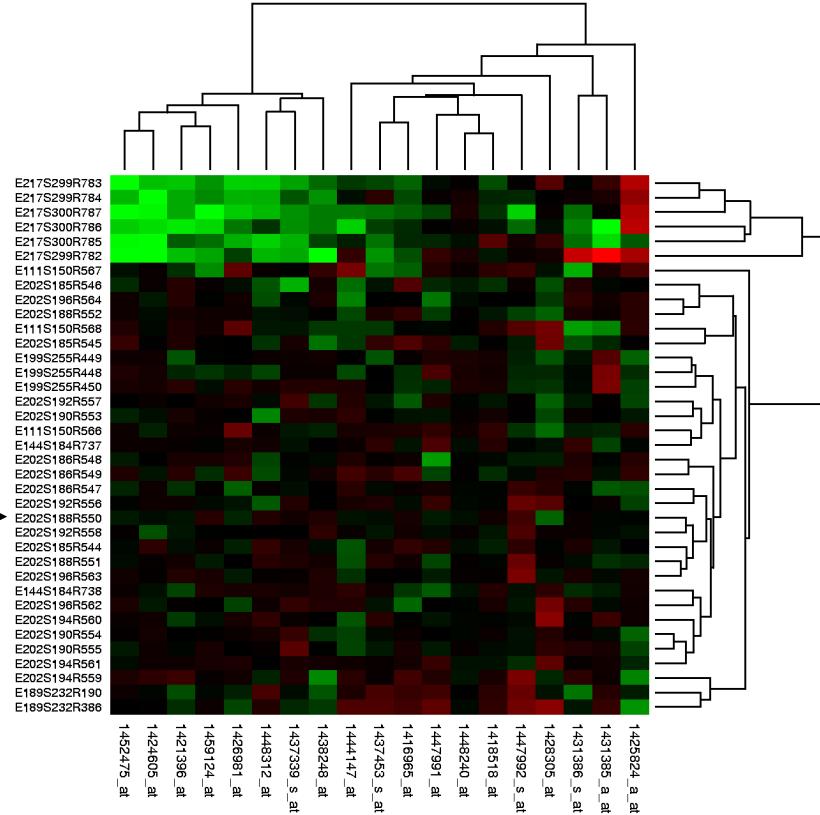
(insane, but feasible with a big enough blackboard)

Basic Local Alignment and Search Tool

(I bet it was a BLAST after trying to do things by hand)

“Heat” maps

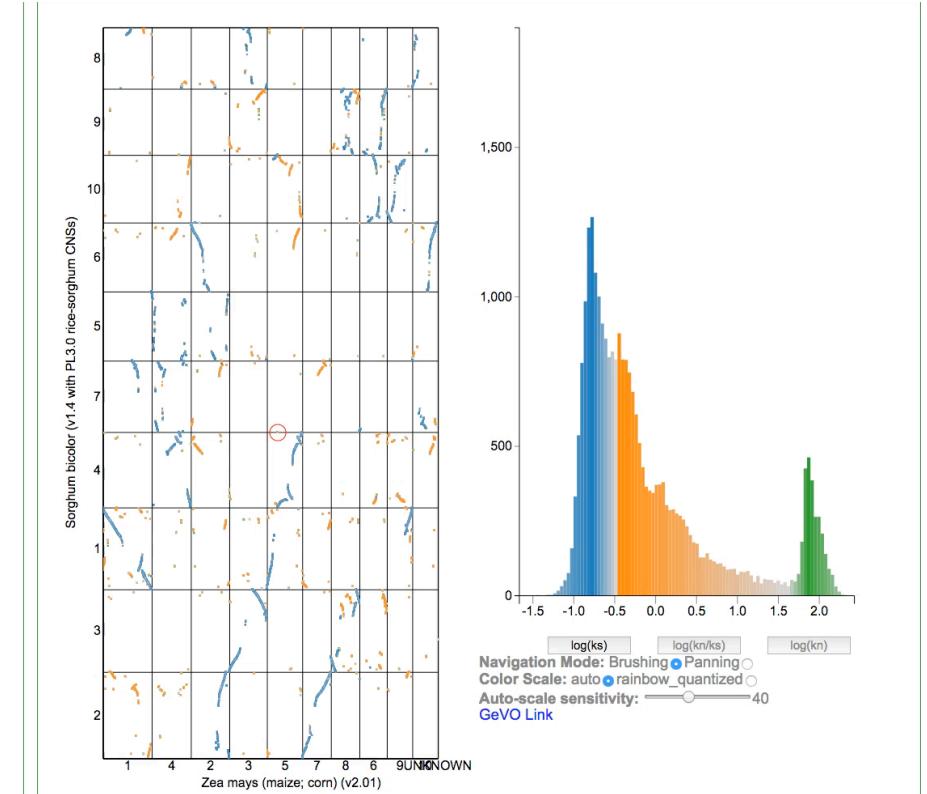
(Red and green? What the \$#%@!?)



Example #1: Synteny Plots

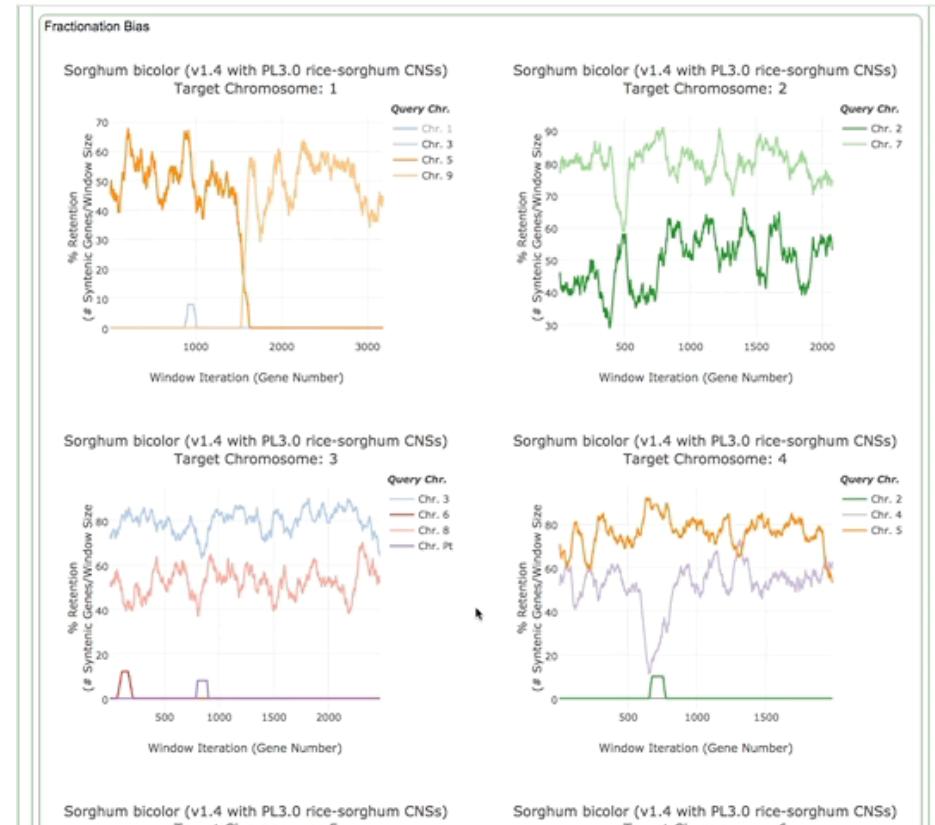
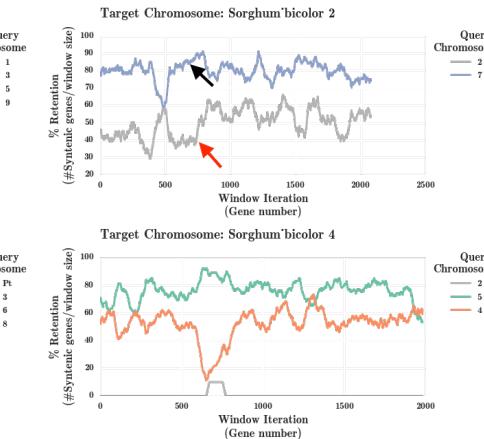
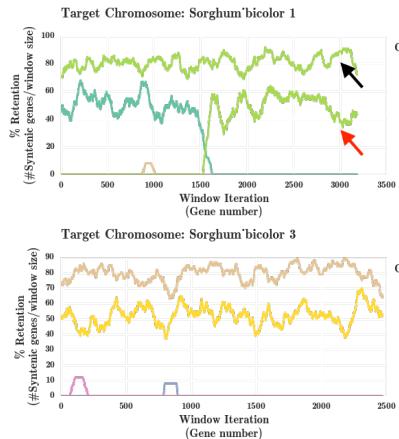


* NOT Plotly!



@seanastephens
<https://github.com/hdc-arizona/synteny-vis>

Example #2: Fractionation Bias



Thanks, Plotly!

Visions of the Future (aka Asher and Blake's wish list)

Don't forget all that chaos!

- 1) **Multidimensional analysis (N-genome, N-dataset comparisons)**
- 2) Parallelization of everything for more computational power. Literally everything
- 3) **Build systems that move, analyze, and visualize data + interoperate**
- 4) Death to single-purpose analysis tools
- 5) Online manuscripts, including data hosting & new (interactive) visualizations
- 6) Convince Microsoft to remove visualizations from Excel
- 7) Rule the world with an iron fist

Case Study #1

Assessing the Age of Duplicated Genes



Synonymous -vs- Non-Synonymous Mutation

...a simplified example using my favorite gene...

RNA!
"transcript"

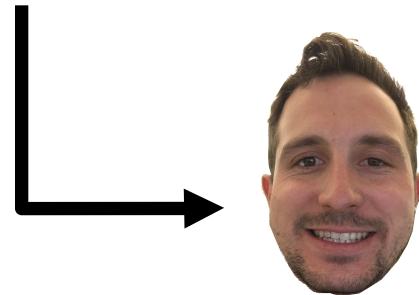
AUGGCUUCUCAUGAACGUUAG

start!

A - S - H - E - R

stop!

Protein!



		Second Position							
		U	C	A	G				
		code	Amno Acid	code	Amno Acid	code	Amno Acid	code	Amno Acid
U	UUU	phe	UCU	tyr	UGU	cys	U		
	UUC		UCC		UAC		C		
	UUA	leu	UCA	ser	UAA	STOP	A		
	UUG		UCG		UAG	STOP	G		
C	CUU		CCU		CAU		GGU		
	CUC		CCC		CAC		GGC		
	CUA		CCA		CAA		CGA		
	CUG		CCG		CAG		CGG		
A	AUU		ACU		AAU		AGU		
	AUC		ACC		AAC		AGC		
	AUA	ile	ACA	thr	AAG	lys	AGA	ser	U
	AUA		ACG		AAA		AGG		C
G	AUG	met	ACG		AAG		AGG		A
	GUU		GCU		GAU		GGU		G
	GUU		GCC		GAC		GGC		
	GUU		GCA		GAA		GGA		
T	GUG		GCG		GAG		GGG		

Synonymous Mutation (Ks)

...protein product **doesn't change...**



AUGGCUUCCUCAUGAACGGUAG

start! A S H E R *stop!*



		Second Position													
		U			C			A			G				
		code	Amino Acid	code	Amino Acid	code	Amino Acid	code	Amino Acid	code	Amino Acid	code	Amino Acid		
U	UUU	phe	UCU	tyr	UGU	cys	U								
	UUC		UCC		UAC		C								
	UUA		UCA	ser	UAA	STOP	UGA	STOP							
	UUG	leu	UCG	UAG	STOP	UGG	tpp								
C	CUU		CCU		CAU		CGU								
	CUC		CCC		CAC		CGC								
	CUA		CCA	pro	CAA	gln	CGA								
	CUG		CCG		CAG		CGG								
A	AUU		ACU		AAU		AGU								
	AUC		ACC		AAC		AGC								
	AUA		ACA	ile	AAT	asn	AGA	ser							
	AUG	met	ACG		AAA		AGG								
G	GUU		GCU		GAU		GGU								
	GUU		GCC		GAC		GGC								
	GUU		GCA		GAA		GGG								
	GUU		GCG	val	GAG	glu	GGG	gly							

Non-Synonymous Mutation (Kn)

...protein product changes...



AUGGCUUCUCAUGAAAGUUAG

start!

A S H E S

stop!



+ time



AUGGCUUCUCACCGAACGCUAG

start!

A S H E R

stop!

$$K_S = 2$$

+ more time



$\kappa_s = 4$

		Second Position					
		U	C	A	G		
		code	Amno Acid	code	Amno Acid	code	Amno Acid
U	UUU	phe	UCU	tyr	UGU	cys	U
	UUC		UCC		UAC	UGC	C
	UUA	leu	UCA	ser	UAA	STOP	A
	UUG		UCG		UAG	STOP	G
C	CUU		CCU	CAU	CGU		U
	CUC		CCC	CAC	CGC		C
	CUA	leu	CCA	CAA	CGA	arg	A
	CUG		CCG	CAG	CGG		G
A	AUU		ACU	AAU	AGU		U
	AUC	ile	ACC	AAC	AGC	ser	C
	AUA		ACA	AAA	AGA	A	A
	AUG	met	ACG	AAG	AGG	G	G
G	GUU		GCU	GAU	GGU		U
	GUU		GCC	GAC	GGC		C
	GUU		GCA	GAA	GGA	A	A
	GUG		GCG	GAG	GGG		G

WHAT DOES IT MEAN THOUGH?!?!

Synonymous mutations do not affect fitness (no functional changes)

...thus, synonymous mutations are invisible to selection.

...thus, synonymous mutations can accumulate at a roughly linear rate.

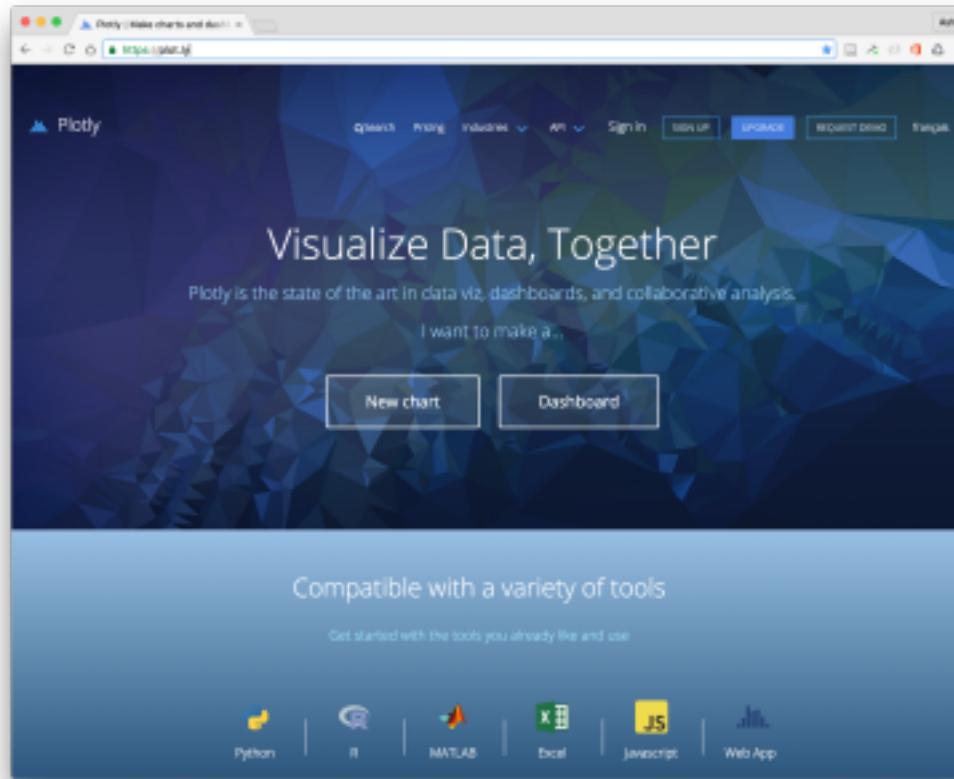
...thus, **more synonymous mutations = more time since gene duplication**



$K_s \sim$ Molecular Clock



Getting Started with Plot.ly

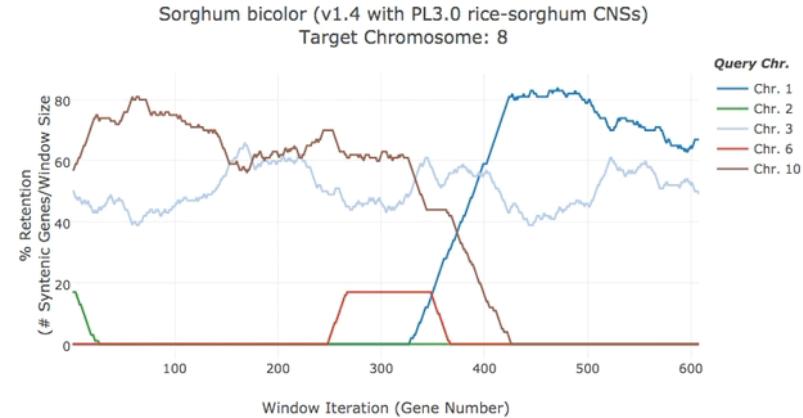




What IS Plot.ly?

Plot.ly is an open-source* data visualization tool.

- ... was built using **Python** (Django) & JavaScript.
- ... offers a web application for **visualization & analysis**.
- ... provides plotting **APIs** for many popular languages.
- ... plots are **fully interactive**, and rendered with D3.js or WebGL (for 3D).
- ... free, paid, and on-site offerings.



* well, *mostly* open-source.

Step 1: *Install Plotly API Package*



Install with PIP!

```
(sudo) pip install plotly
```

Update with PIP!

```
(sudo) pip install plotly --upgrade
```



Install with CRAN!

```
install.packages("plotly")
```

Or, with **DevTools & GitHub!**

```
devtools::install_github("ropensci/plotly")
```

Step 2: Setup your API keys

...these are necessary for communicating with Plot.ly's servers.

Question: *What is an API?*

Answer: *An “Application Programming Interface”*

Question: ...?

Slightly More Useful Answer: *A way to interact with someone else’s program, using programming languages rather than a graphical user interface*



```
import plotly  
plotly.tools.set_credentials_file( \  
    username='YourAccount', \  
    api_key='YourKey')
```

NOTE: *You only need to do this once on your computer account, even if you use different virtual environments.*



```
Sys.setenv("plotly_username"="user")  
Sys.setenv("plotly_api_key"="pass")
```

NOTE: *I don't know if you need to do this every session or only once.*



Step 3: Compose a Plot

Figure() - The Plotly “Plot” Object

Composed of two parts: ‘Data’ and ‘Layout’

Data()

Contains the information to be plotted.

Composed of ‘trace’s

Layout()

Contains information about the plot

i.e. title, labels, fonts, annotations, etc.

Figure()

Combines Data() and Layout()

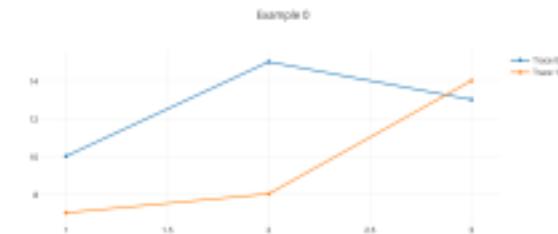
```
import plotly.plotly as py
from plotly.graph_objs import *

trace0 = Scatter(x=[1, 2, 3],y=[10, 15, 13])
trace1 = Scatter(x=[1, 2, 3],y=[7, 8, 14])
my_data = Data([trace0, trace1])

my_layout = Layout(
    title="Example 0",
    xaxis=dict(title="X"))

my_figure = Figure(
    data=my_data, layout=my_layout)

py.plot(my_figure, filename='example0')
```



<https://plot.ly/~asherkhb/462/example-0/>



Step 3: Compose a Plot

plotly charts are described *declaratively*...

1. plotly::plot_ly
2. plotly::add_trace
3. plotly::layout

Every aspect of a chart (the colors, the grid-lines, the data, and so on) has a corresponding key in these call signatures.



<https://plot.ly/r/reference/>

Step 4: Plot!!



```
# Import Online Module  
import plotly.plotly as py  
  
# Import Offline Module  
Import plotly.offline as pyo
```

Regular Python

```
py.plot()    # For ONLINE  
pyo.plot()   # For OFFLINE
```

iPython

```
py.iplot()    # ONLINE (still inline)  
pyo.iplot()   # OFFLINE (inline)
```



Default is offline, simply print!
Interactive plot will open either in web browser
or in Rstudio viewer

```
>>> library(plotly)  
>>> p <- plot_ly(midwest,  
                      x=~percollege,  
                      color = ~state,  
                      type = "box")  
>>> p
```

To save your plots online to Plotly's servers...

```
>>> plotly_POST(p, filename = "r-  
docs/midwest-boxplots")
```

Online vs Offline Plotting

Online Plotting

- Data is sent to Plotly's servers, plots are generated and links are returned.
- Embedded plots can be returned to iPython.

Pros

- Ability to adjust look of plots using GUI.
- Easy to share plots with other users.

Cons

- Internet connection required, need to set up API keys.
- Data is public (unless Pro account).

Offline Plotting

- Data is bundled with Plotly's library locally, made available as local .html files or divs.
- Embedded plots are available when using iPython.

Pros

- No internet connection or keys required.
- Data remains private without paying.

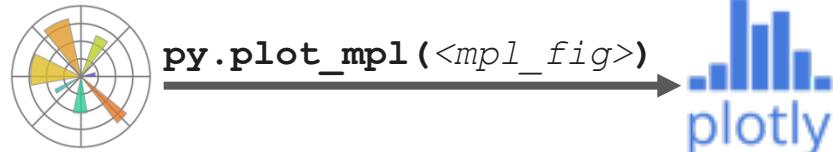
Cons

- Must programmatically adjust all visual parameters to achieve desired looks
- Some of Plotly's functionality not available.

Other Features & Benefits of Plotly

→ Matplotlib converter!

<https://plot.ly/matplotlib/>



→ Live & static dashboards!!

<https://plot.ly/dashboards/>



→ Database connections!*

<http://help.plot.ly/database-connectors/>

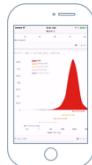
→ Presentations?!*

<https://formidable.com/open-source/spectacle-editor/>

A promotional image for Spectacle Editor. At the top, it shows the logos for Formidable and Plotly, followed by the text "PRESENT". Below this is the Spectacle Editor logo, which features a blue bar chart icon. The text "Spectacle Editor" is displayed above the tagline "The world's first open-source presentations editor." To the right of the main title, there is a list of database connection options, each with a small icon and text:

- SQLite: Connecting to a SQLite Database
- Amazon Redshift: Connecting to a Redshift Database
- PostgreSQL: Connecting to a Postgres Database
- MySQL: Connecting to a MySQL Database
- Microsoft SQL Server: Connecting to a MS SQL Server Database
- MariaDB: Connecting to a MARIA DB Database

→ Mobile Ready!



* untested by me personally.

Getting more **help**...

Check out the docs for basics & examples...

Python: <https://plot.ly/python/>

R: <https://plot.ly/r/>

If you can't figure it out, post to the forums...

<http://community.plot.ly/>