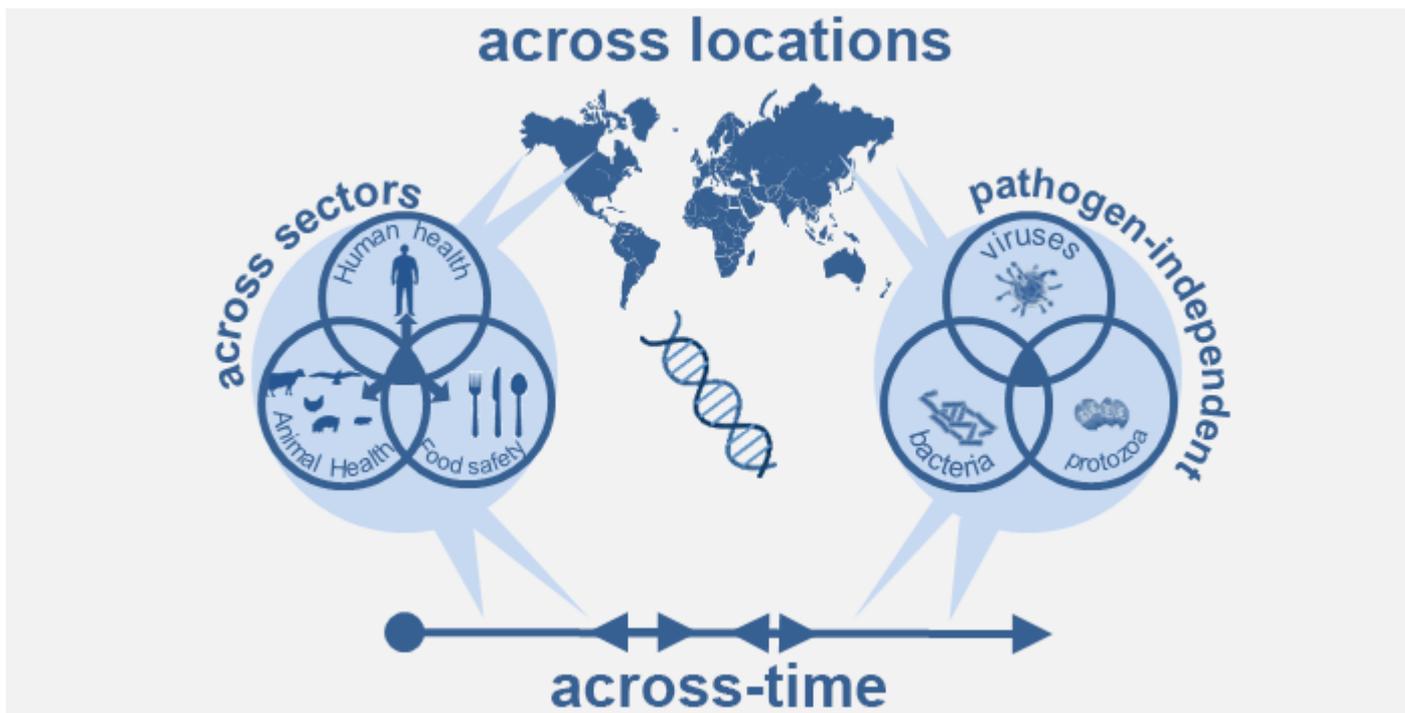
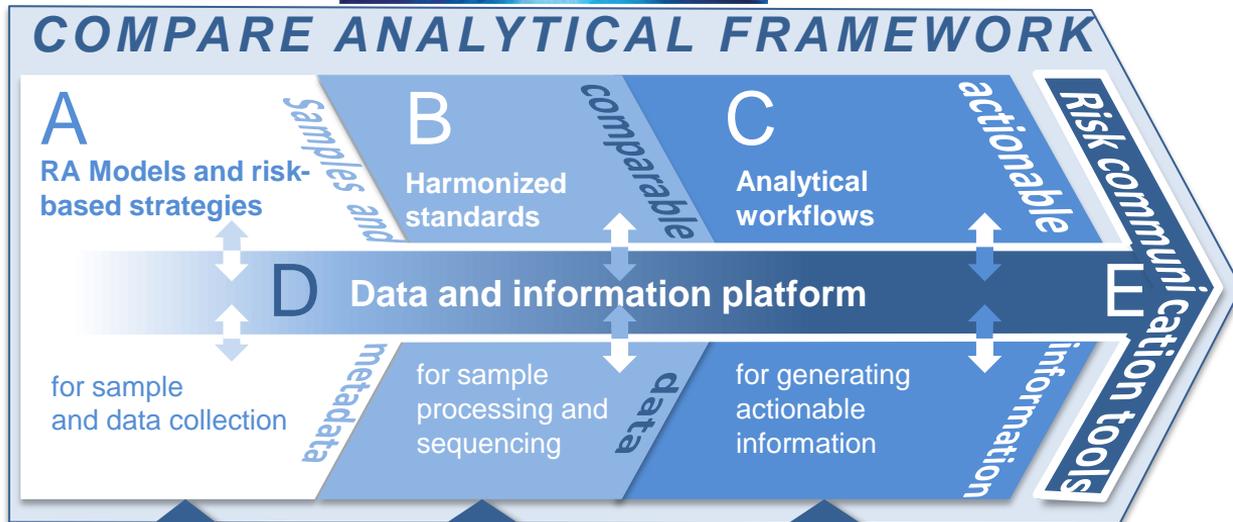


COLLABORATIVE PLATFORM FOR GENOMICS BIG DATA ANALYSIS @ MTA CLOUD

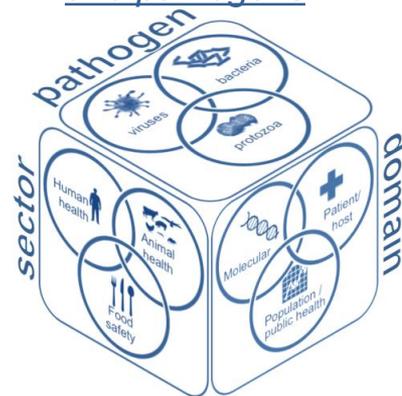
ISTVAN CSABAI, LASZLO OROSZLANY, JANOS SZALAI-GINDL, DAVID VISONTAI, LASZLO DOBOS, DEZSO RIBLI
EÖTVÖS UNIVERSITY DEPT. OF PHYSICS OF COMPLEX SYSTEMS & WIGNER RCP

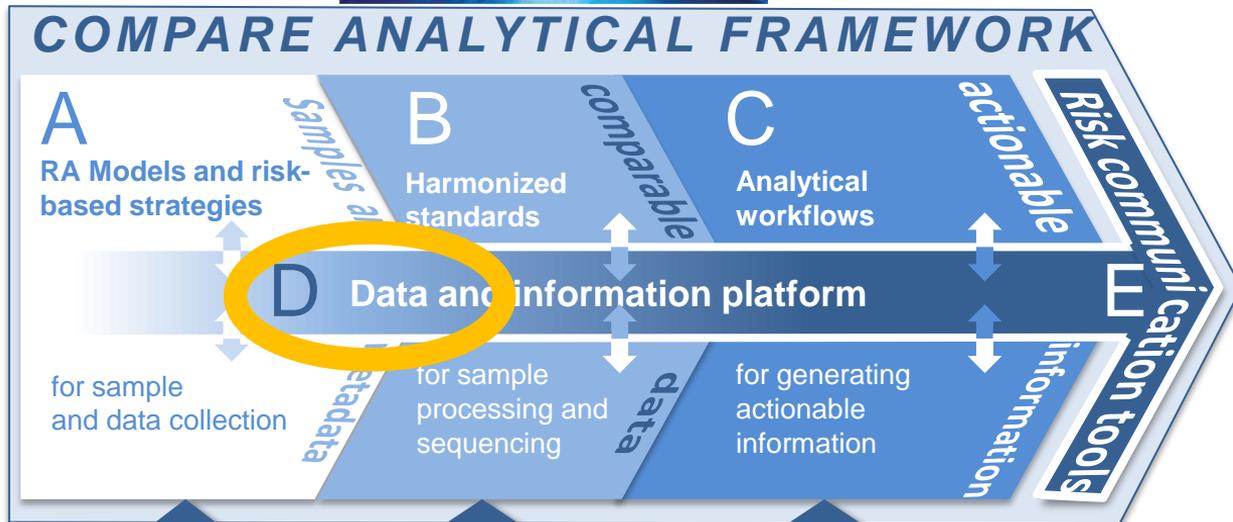




- User-stakeholders:
- Public and Veterinary Health, and Food Safety Authorities, Wildlife/Nature management;
 - Clinicians, General Practitioners, Veterinary practices;
 - Scientific community, Research Labs (Microbiological, Clinical, Food safety labs);
 - Food industry;

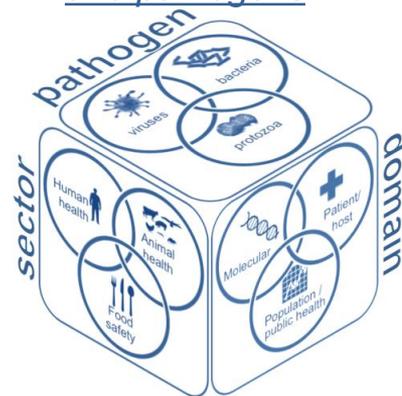
Across domains, sectors and pathogens



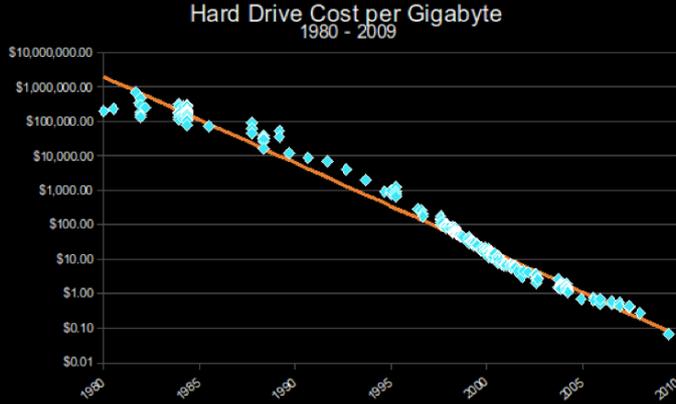


- User-stakeholders:
- Public and Veterinary Health, and Food Safety Authorities, Wildlife/Nature management;
 - Clinicians, General Practitioners, Veterinary practices;
 - Scientific community, Research Labs (Microbiological, Clinical, Food safety labs);
 - Food industry;

Across domains, sectors and pathogens

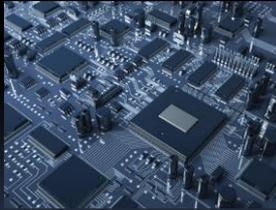


Exponentially cheaper devices



Moore's law

Exponential growth



Elektronics

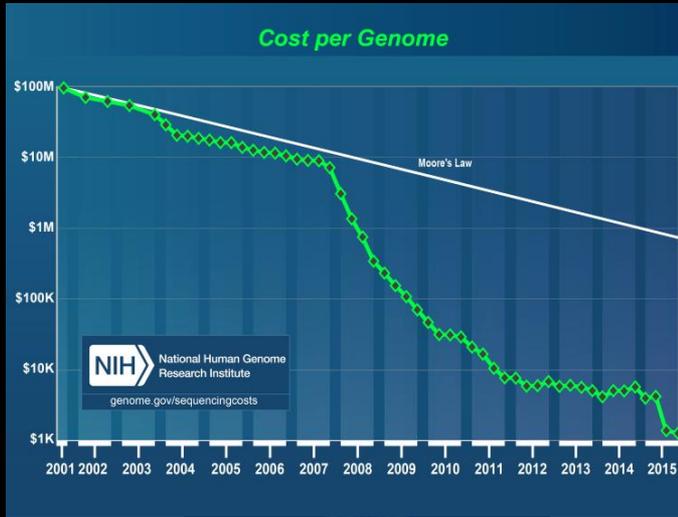


Sensors



Data

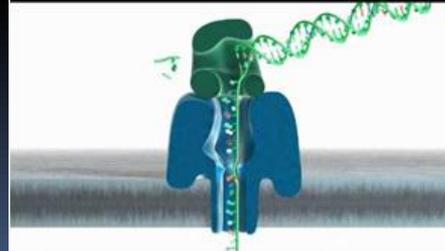
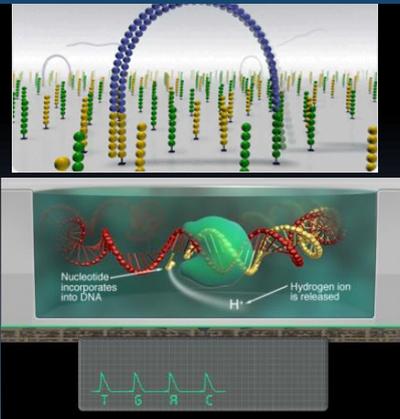
Moore's law in genetics



Human genome sequencing
1990-2003: 13yrs /2.7 Bn USD
2016: ~days/1000 USD
2020: ??????

CCD!

- 2006 X Prize 10M, 100 genomes
30 days, \$10k - cancelled
- Microarray
- Mass spectroscopy
- Digital microscopy
- ...



Oxford Nanopore
100Mb, \$900



The Cancer Genome Atlas Data: Navigating the Data Portal and the Cancer Genomics Hub

The Cancer Genome Atlas
<http://cancergenome.nih.gov/>

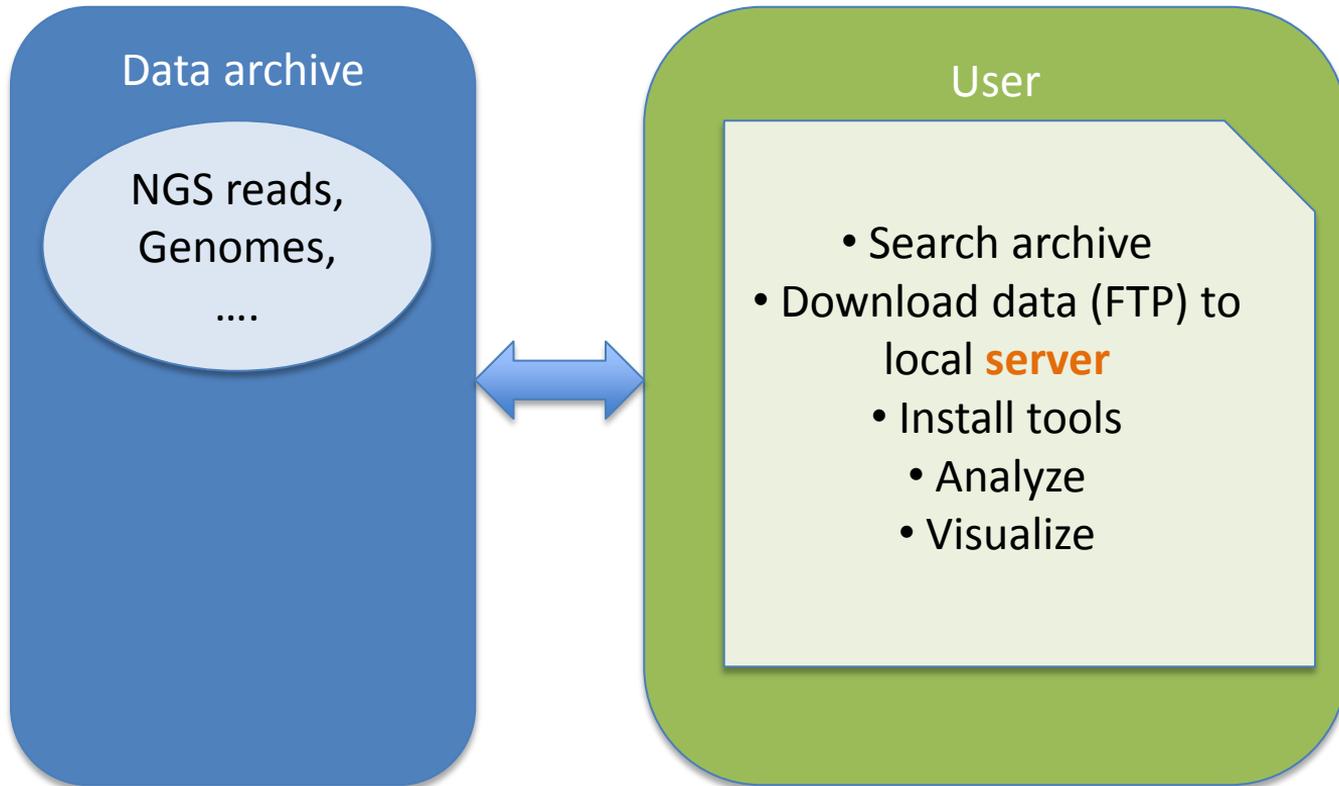
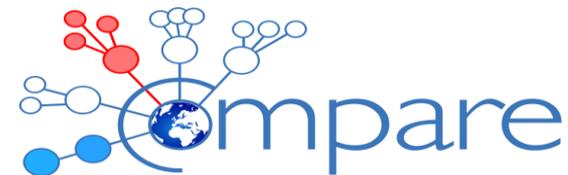
3.2Bn nucleotides / human genome

The Cancer Genome Atlas (TCGA) is a large-scale collaborative effort led by the National Institutes of Health to map the genomic changes that occur in over 30 types of human cancer, including nine rare tumors. Its goal is to support new discoveries and accelerate the pace of research aimed at improving the diagnosis, treatment, and prevention of cancer.

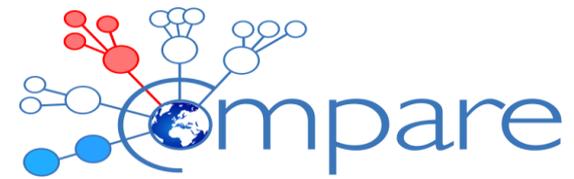
TCGA is a community resource project. The information generated by TCGA is centrally managed and entered into databases as it becomes available, making the data rapidly accessible to the entire research community. By January 2014, TCGA had generated one petabyte of data for about 10,000 cases of tumor and matching normal tissue samples.

TCGA data are available in two data repositories: the TCGA Data Portal and the Cancer Genomics Hub. All data can be accessed directly from the TCGA Data Portal regardless of which repository houses the data file.

Data exploration: Traditional approach



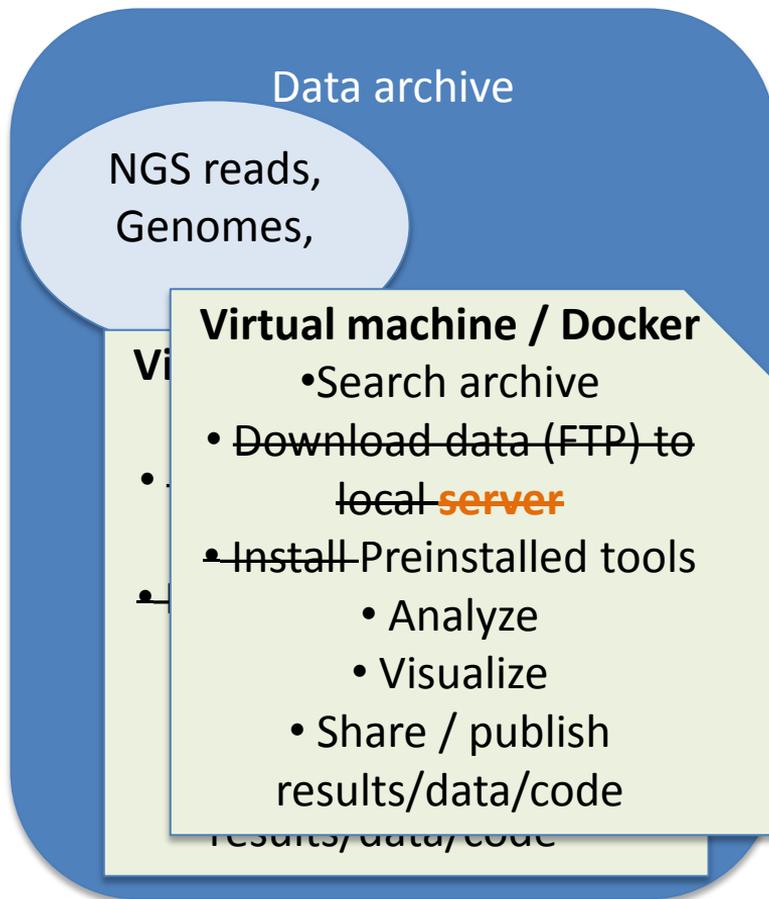
Challenges



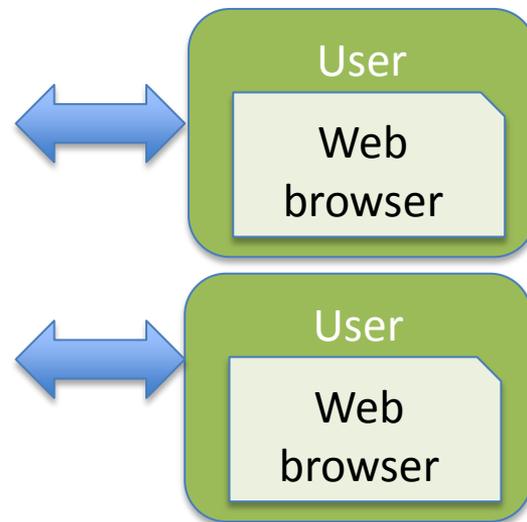
- Big Data – downloading data is not optimal/possible
- Data sharing is not enough – share data + complete processing pipeline + result figures, tables, ...



Cloud+notebook approach



Virtualization: OpenStack + Docker
“Workflow”: Jupyter/Ipython notebook



Kooplex

Infrastructure for flexible collaboration



Worksheets

Run pre-compiled worksheets to hide program code and focus on the problem.

Worksheets hide the complexity of notebooks from end users interested in scientific output.

[→ to worksheets](#)



Jupyter notebooks

Run existing or create Jupyter notebooks from existing projects to process your data or author your own projects, notebooks and share with others.

[→ Browse notebooks](#)



Gitlab

Manage your project, add members to it, file issues.

[→ to Gitlab](#)



Owncloud

Manage easily your data files through the owncloud... (just as easy as in Dropbox! :)

[→ to Owncloud](#)

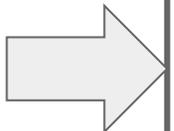
-
-

[About](#)
[Contact](#)

Designed and built with all the love in the world by the **Kooplex Team**. We try to maintain it too! :)

ARCHITECTURE

INTERNET



NGINX

ownCloud
NFS
LDAP
GitLab

HUB

C-HTTP-P

Jupyter
NOTEBOOKS

Jupyter
DASHBOARD

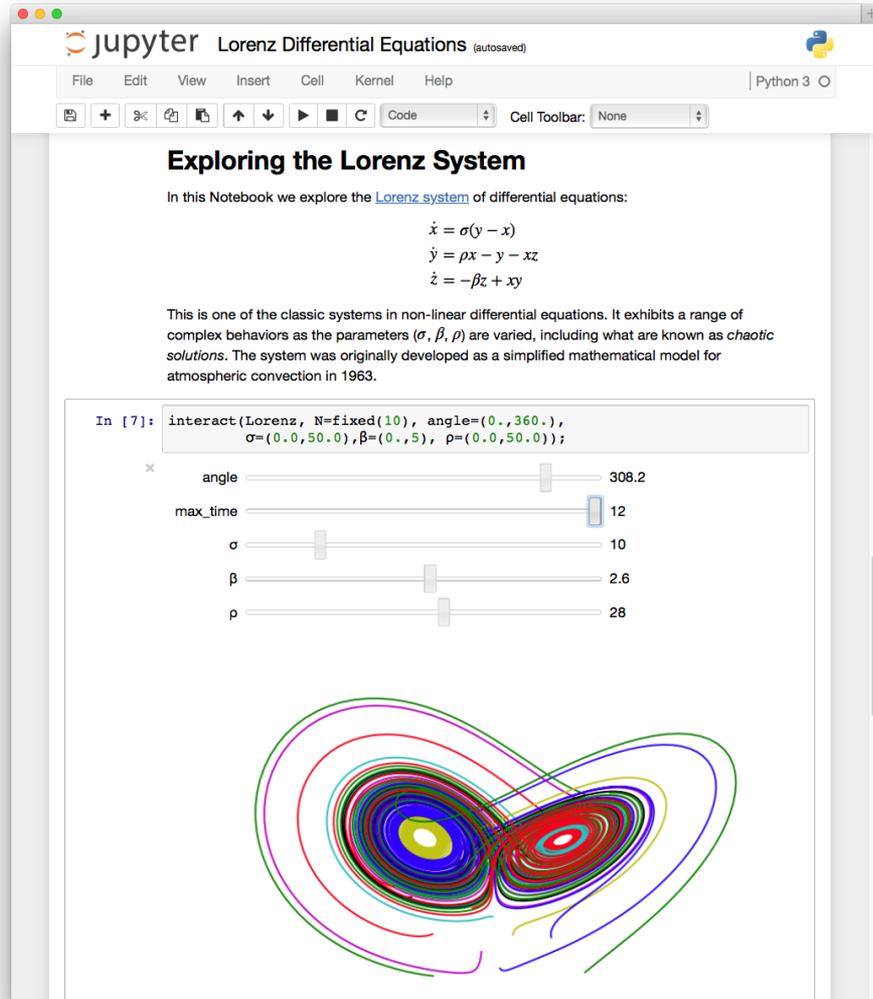
Compute capability
/
Databases

COMPONENTS

- Nginx/C-HTTP-P [networking/routing, user↔resources]
- Docker [compartmentalization of services and notebooks]
- LDAP [authentication and user database]
- NFS [persistent user workspace]
- HUB [front page + service spawner, django powered website]
- Gitlab [project management, versioning, collaborative forum]
- Jupyter [dev. env., sandbox, DASHBOARDS]
- Own Cloud [small data drag & drop, sharing]
- ... + [workflow/pipeline services]

JUPYTER

- Integrate:
 - Executable code (Python, R, Matlab, SQL, bash, ...)
 - Notes, descriptions, math equations, ...
 - Results, figures, interactions
- Accessible from anywhere through browser
- Runs in the cloud



The screenshot shows a Jupyter Notebook interface with the title "Lorenz Differential Equations (autosaved)". The notebook content includes a title "Exploring the Lorenz System", a brief introduction, the Lorenz equations, a description of the system, and an interactive plot of the Lorenz attractor.

Exploring the Lorenz System

In this Notebook we explore the [Lorenz system](#) of differential equations:

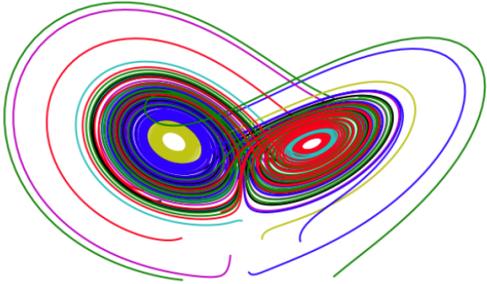
$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= \rho x - y - xz \\ \dot{z} &= -\beta z + xy\end{aligned}$$

This is one of the classic systems in non-linear differential equations. It exhibits a range of complex behaviors as the parameters (σ , β , ρ) are varied, including what are known as *chaotic solutions*. The system was originally developed as a simplified mathematical model for atmospheric convection in 1963.

```
In [7]: interact(Lorenz, N=fixed(10), angle=(0., 360.),
                sigma=(0.0, 50.0), beta=(0., 5), rho=(0.0, 50.0));
```

The interactive plot shows the Lorenz attractor with the following parameters:

Parameter	Value
angle	308.2
max_time	12
σ	10
β	2.6
ρ	28



EBOLA

www.compare-europe.hu/ena_europe_loco.html

```
In [12]: width, height = 650, 500
flu_map = folium.Map(location=[47, -17], zoom_start=3,
                    tiles='OpenStreetMap', width=width, height=height)
```

Add point to the map object

- Lets make point area proportional to number of cases
 - This is misleading, because somewhere all the cases around have the same position (Europe), and somewhere the positions are more scattered (Shanghai)

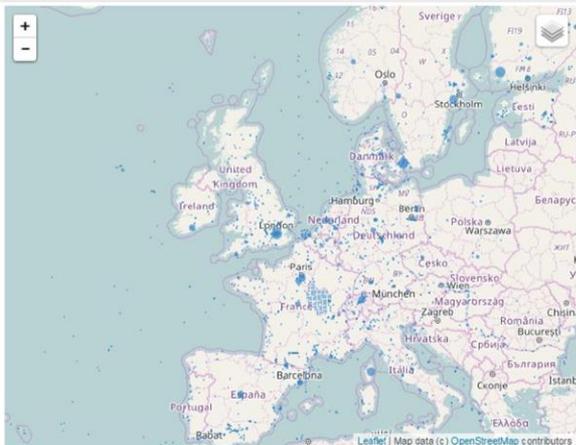
```
In [13]: for i in xrange(len(uniq_locs_w_acc)):
loc=(uniq_locs_w_acc.iloc[i]['lat'],uniq_locs_w_acc.iloc[i]['lon'] )
name='Number of cases: '+str(uniq_locs_w_acc.iloc[i]['count'])
name+=' Accessions: '+uniq_locs_w_acc.iloc[i]['acc_list']
size=uniq_locs_w_acc.iloc[i]['count'] ** 0.5

flu_map.circle_marker(location=loc, radius=1e3*size,
                    line_color='none',fill_color='#3186cc',
                    fill_opacity=0.7, popup=name)
```

And finally draw the map

```
In [15]: inline_map(flu_map)
```

Out[15]:

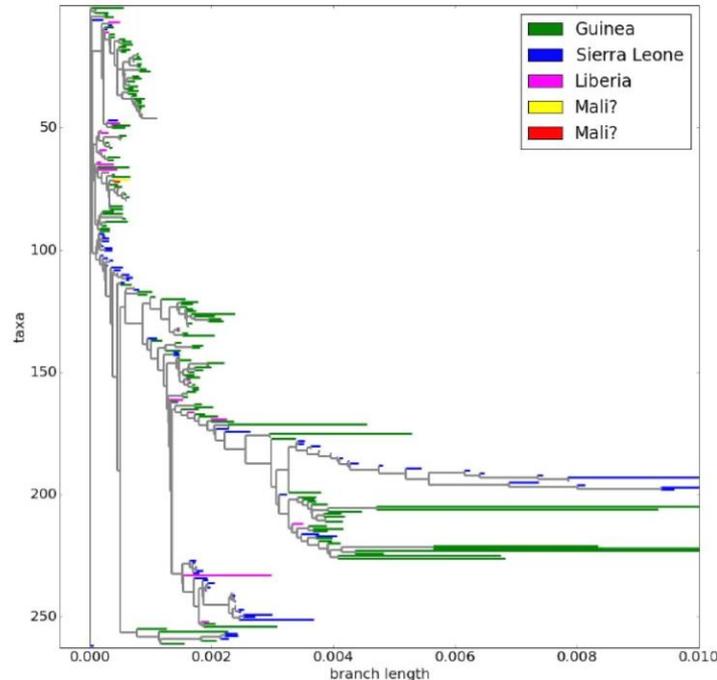


```
!matplotlib inline

#some settings
matplotlib.rcParams['font', size=20]
matplotlib.rcParams['lines.linewidth'] = 3
matplotlib.rcParams['figure.figsize'] = (16,16)

fig,ax=plt.subplots()
#custom Legend
gui_proxy = plt.Rectangle((0, 0), 1, 1, fc="green")
sie_proxy = plt.Rectangle((0, 0), 1, 1, fc="blue")
lib_proxy = plt.Rectangle((0, 0), 1, 1, fc="magenta")
dpr1_proxy = plt.Rectangle((0, 0), 1, 1, fc="yellow")
dpr1_proxy = plt.Rectangle((0, 0), 1, 1, fc="red")
ax.legend([gui_proxy,sie_proxy,lib_proxy,dpr1_proxy],
          ['Guinea', 'Sierra Leone', 'Liberia', 'Mali?'])

#draw tree
def my_label(clade):
    return None
Phylo.draw(tree,my_label,axes=ax,xlim=(-0.0005,0.01))
```



GLOBAL SEWAGE METAGENOME

 jupyter sewage-metadata Last Checkpoint: 16 minutes ago (autosaved)



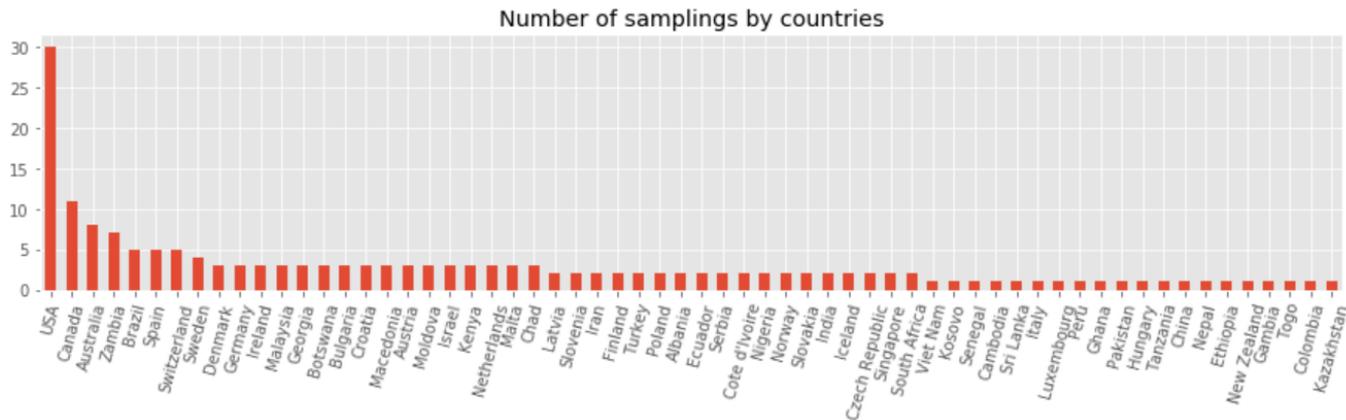
File Edit View Insert Cell Kernel Help

Python 3

          Markdown   CellToolbar

Number of samplings

```
In [4]: sewage_metadata['country'].value_counts().plot.bar(figsize=(15,3), title="Number of samplings by countries")
plt.xticks(rotation=75)
plt.show()
```



```
In [5]: # For label purposes
sewage_metadata['date_fact'] = sewage_metadata['collection_date']
sewage_metadata['type_fact'] = sewage_metadata['sewage_type']
sewage_metadata = sewage_metadata.apply(lambda x: (pd.factorize(x)[0]) if x.name in ['date_fact', 'type_fact'] else x)
```

COSMOLOGY WITHOUT DARK ENERGY

jupyter darkEnergy Last Checkpoint: 3 minutes ago (unsaved changes)

File Edit View Insert Cell Kernel Help

⏏ + ⏪ ⏩ ⏴ ⏵ ⏴ ⏵ Code CellToolbar

Calculate ISW for non-ΛCDM model without Dark Energy

Growth function

The definition of growth function:

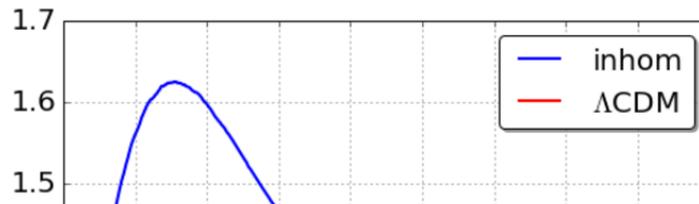
$$D_1(z) = \frac{H(z)}{H(0)} \int_z^\infty \frac{(1+z)}{H^3(z)} dz' \left[\int_0^\infty \frac{(1+z)}{H^3(z)} dz' \right]^{-1}$$

EdS:

$$\int_{9.0}^\infty \frac{(1+z)}{H^3(z)} dz' = \frac{27}{8} \int_{9.0}^\infty \frac{(1+z)}{(1+z)^{9/2}} dz' = \frac{27}{8} \left[\frac{-2}{5(1+x)^{5/2}} \right]_{9.0}^\infty = \frac{27}{8} \frac{1}{250\sqrt{10}} = 0.0042690375$$

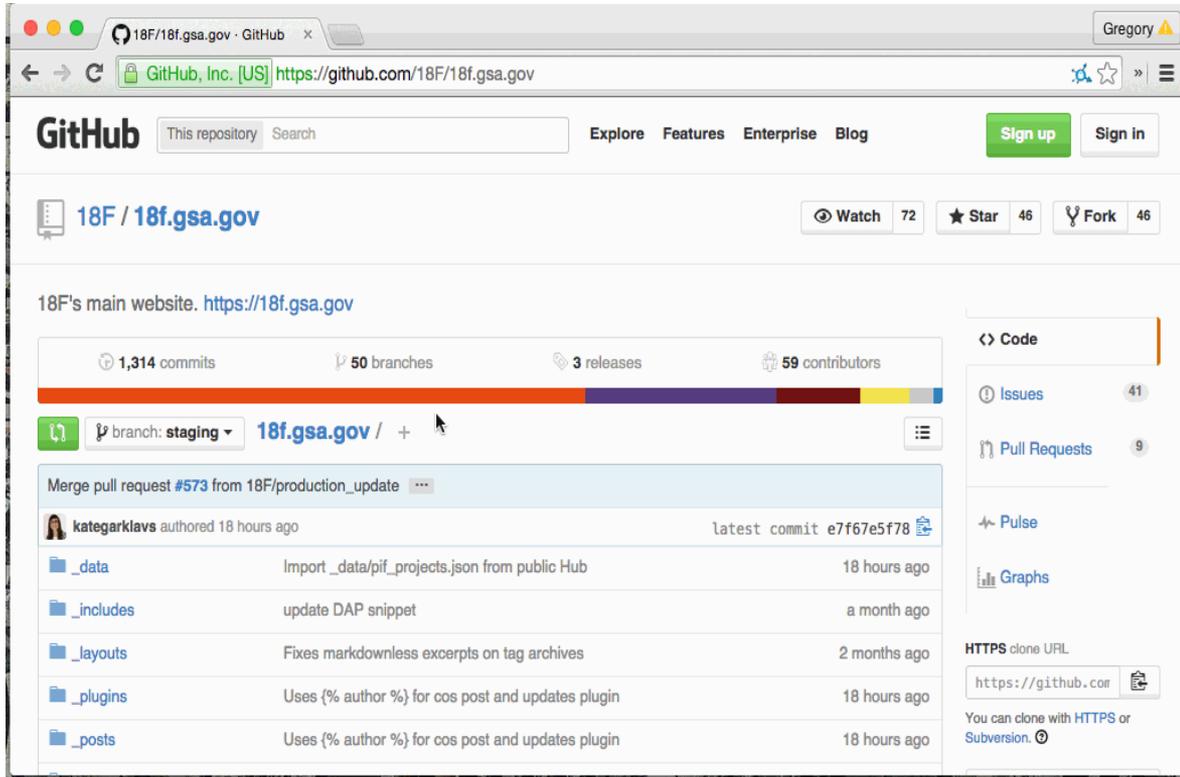
$$\int_{0.0}^\infty \frac{(1+z)}{H^3(z)} dz' = 1.35$$

In [36]: `plotModel(x)`



GITHUB/GITLAB

- project management, versioning,
- collaborative forum



The screenshot shows a web browser displaying the GitHub repository page for `18F/18f.gsa.gov`. The browser's address bar shows the URL `https://github.com/18F/18f.gsa.gov`. The page header includes the GitHub logo, a search bar, and navigation links for `Explore`, `Features`, `Enterprise`, and `Blog`. There are `Sign up` and `Sign in` buttons in the top right corner. Below the header, the repository name `18F / 18f.gsa.gov` is displayed, along with statistics: `72` Watch, `46` Star, and `46` Fork. A link to the main website is provided: `18F's main website. https://18f.gsa.gov`. A progress bar shows repository statistics: `1,314` commits, `50` branches, `3` releases, and `59` contributors. Below the progress bar, there is a branch selector set to `branch: staging` and a `+` button to create a new branch. A merge pull request `#573` from `18F/production_update` is highlighted. The commit history is shown with the following entries:

Commit	Message	Time
<code>kategarklavs</code>	latest commit <code>e7f67e5f78</code>	18 hours ago
<code>_data</code>	Import <code>_data/pif_projects.json</code> from public Hub	18 hours ago
<code>_includes</code>	update DAP snippet	a month ago
<code>_layouts</code>	Fixes markdownless excerpts on tag archives	2 months ago
<code>_plugins</code>	Uses <code>{% author %}</code> for cos post and updates plugin	18 hours ago
<code>_posts</code>	Uses <code>{% author %}</code> for cos post and updates plugin	18 hours ago

On the right side, there are links for `Code`, `Issues` (41), `Pull Requests` (9), `Pulse`, and `Graphs`. At the bottom right, there is a section for `HTTPS clone URL` with the URL `https://github.com` and a note: `You can clone with HTTPS or Subversion.`

DROPBOX/OWNCLOUD

The screenshot shows a mobile browser interface with the URL www.dropbox.com/home. The browser's address bar and tabs are visible at the top. Below the browser, the Dropbox website is displayed. On the left, there is a navigation menu with icons for Dropbox, Sharing, Links, Events, Get Started, and Photos. The main content area features the Dropbox logo and a search bar. Below this is a table listing files and folders. The table has three columns: Name, Kind, and Modified. The files listed include folders like 'academic' and 'bahamas', and various document files such as 'iPhones for trip home.doc', '#card 2013.pdf', 'A New Life.txt', 'aaaa.bak', 'aaaa.txt', 'addresses.otl', 'Authors, topics, emails, submissions.pdf', and 'buys.txt'. Each row shows the file name, its kind (e.g., document, file, folder), and the date and time it was last modified.

Name	Kind	Modified
academic	folder	--
iPhones for trip home.doc	document doc	10/27/2011 4:53 AM
#card 2013.pdf	document pdf	12/25/2012 11:44 AM
A New Life.txt	document txt	10/21/2012 4:32 AM
aaaa.bak	file bak	1/29/2013 7:55 AM
aaaa.txt	document txt	2/1/2013 9:34 PM
academic	folder	--
addresses.otl	file otl	3/15/2010 11:09 AM
Authors, topics, emails, submissions.pdf	document pdf	9/15/2009 2:51 AM
bahamas	folder	--
buys.txt	document txt	5/21/2012 5:17 AM



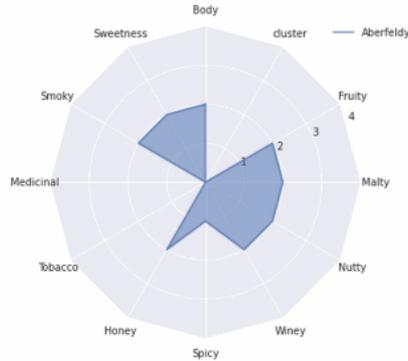
WORKSHEETS/DASHBOARD – ACTIVE SERVICES

Got Scotch?

If you like Aberfeldy you might want to try these five brands. Click one to see how its taste profile compares.

Scotch ▾

	Similarity
Benromach	0.972973
BlairAthol	0.972973
Benrinnes	0.962908
RoyalLochnagar	0.959202
Scapa	0.959043



CRISPR DNA EDITING – GUIDE RNA DESIGNER

Jupyter gRNAdesign (unsaved changes)

File Edit View Insert Cell Kernel Help

Save + Copy Paste Undo Redo Run Stop Refresh Markdown CellToolbar

Start

Motif

Motif: From: Length:

Chr: ChrStart:

Sequence

```
TGAGAGAAGAGTTTCATATTTGCAAGGTCTCAGACATGCCTTTAAAAATTCATAATACTTCTTCTGTGTTTCCCATGCTTGATGG
GGCTCAGAATTGACAGTGACATTTTCAGTAATACAGAGATGAGAAGGGTCAGAGAGGAAGTAAAGTGGCTAGGACTGGCTATG
TAAATTGCCAAGGGGTGCAGTTCAAATGAAAATGTGAGACCTTAATATAAAAATTTTCATGATGGGGATAATGGATCATGAACTA
ATCGTGGGGCTTTGTGCGACACCACGGGACACTTGCCTACAAAGCCATCCCTACTCCCTGCTTAAACGGAGACAAGAACACATT
GGTGATTCTCAAATGCAAACGTCCCTTTACTGAAAACCTAATGAGTTTAAATGCTTACTATACTGTAGTCCCTTATGTGCTGATTA
CATTGTGGAATGCTGCAGGGAGAAAAACAAATTCACCTAATGATGCTAAAGAACCATTAGGAGACTTTACTAGTTTAGGTCACAA
GTGCCTAGAAATACAGAAATGATCTTTGACCTTTCCCTCCATTTTCAAAGGACGCTATTTCTGTGAACCTCATTGACCTCTCCAG
ACAAAGTCCCAGATTGCCTTTGCCAAAGTTAAACCCATCTTGGCTTCTTCCATTTCTCATGTATTCACTAGATCAGTGCATGA
CTGCTGCTCTCAACCCACATCATGGGGCCAGAGCCTTCACTGGTAAATATTTATAAAAAAATACTTTGAAGATTAATCCTTGGTC
AGTAGAGAAAAACTAGACATGGATAGAACAAGAAATGTGGGGTCTGGGCCTCCTCCAGAACTGCCACCACCAGACAATGTT
ATCTTTGACAGATTTGTGGTATCTGGGTGGCTGACTTTTCTTTTGGTGAATAGCAAAAGCCAAAAAAGAGACTGTAACATCT
CATCCATTTTCCCACTTCACTACTCAATTCCTTCCCTTACAGCAAACTATTCAGTCTTTGATCTCTTCCCTTCA
```

