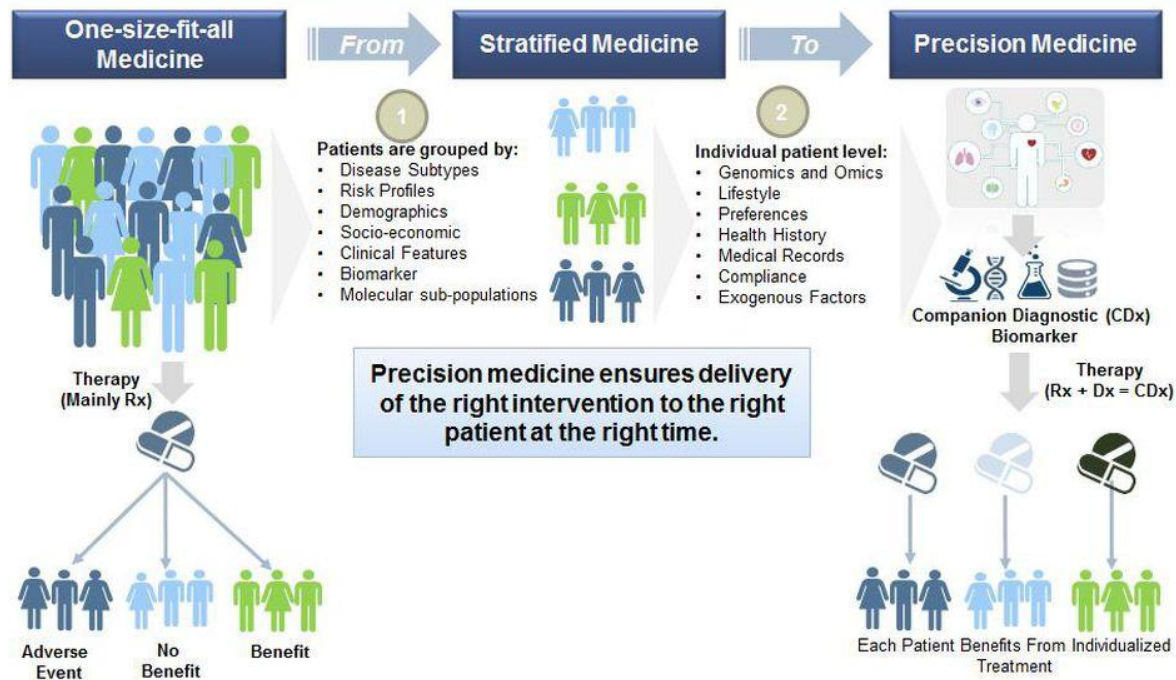The Gen3 data model is flexible and able to host data to be analyzed in different scientific fields.

In this webinar, we will speak about using Gen3 for data analysis in general and show an example of Gen3 used for precision medicine.
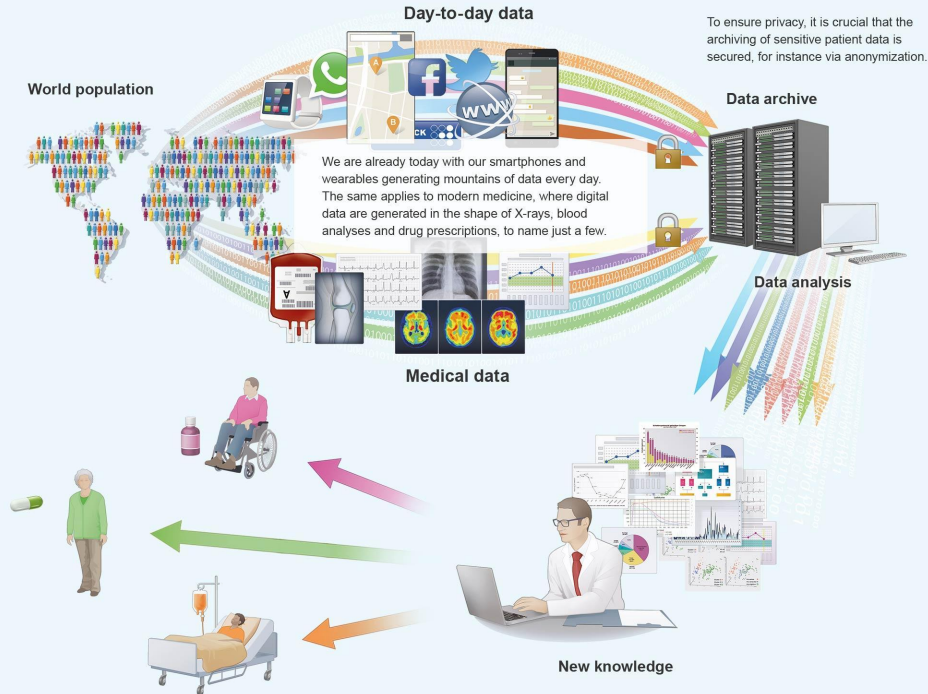
Source: Frost & Sullivan -Figure 1: New Paradigm Shift in Treatment, as referenced in this forbes article

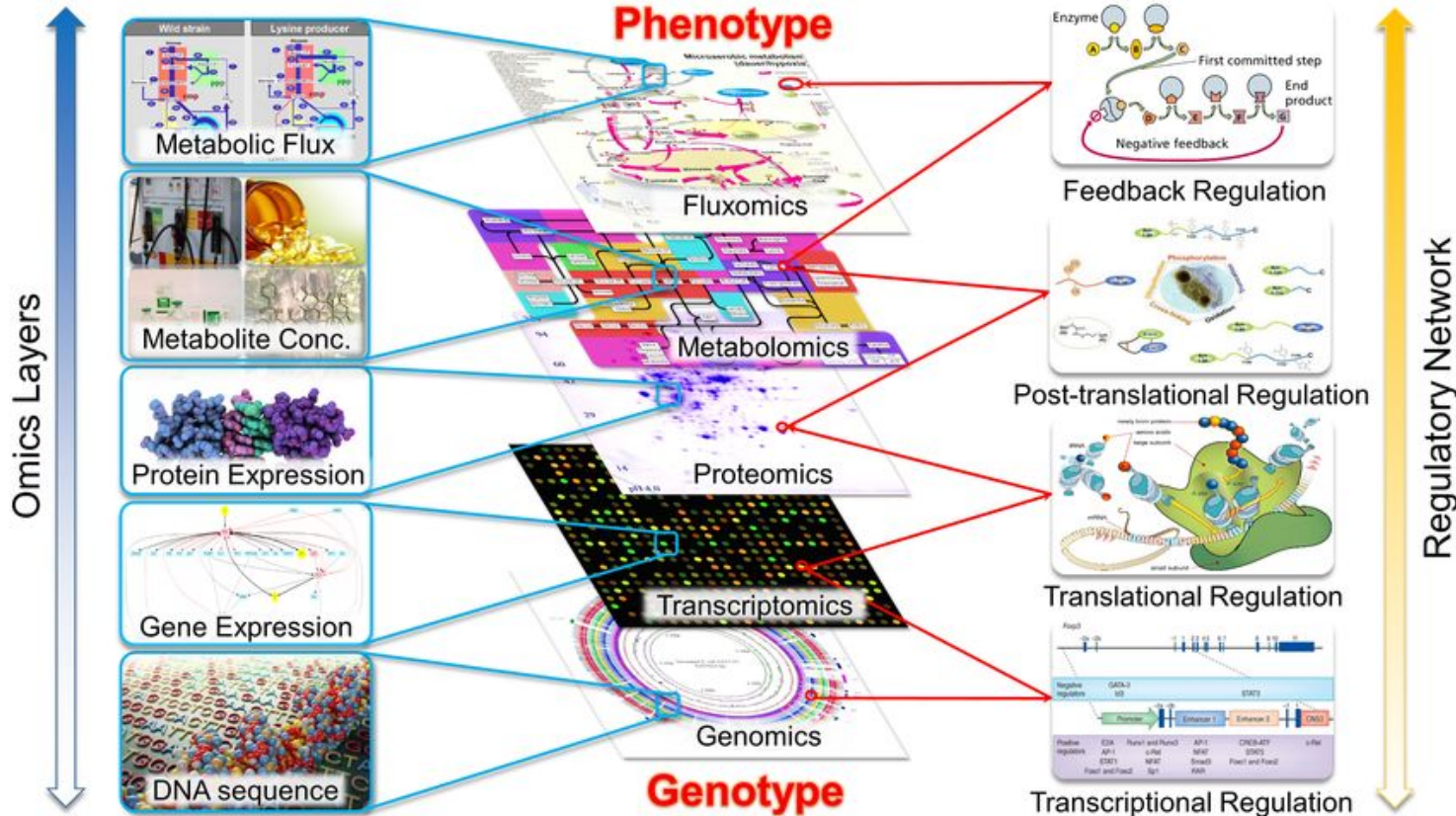Source: Bayer Research, 30 November 2016 Big data in medicine

Source: Guo W, Feng X (2016) OM-FBA: Integrate Transcriptomics Data with Flux Balance Analysis to Decipher the Cell Metabolism.

*Data commons co-locate data, storage and computing infrastructure with commonly used software services, **tools & apps** for analyzing and sharing data to create a resource for the research community.*

Robert L. Grossman, Allison Heath, Mark Murphy, Maria Patterson and Walt Wells, A Case for Data Commons Towards Data Science as a Service, IEEE Computing in Science and Engineer, 2016.   Source of image: The CDIS, GDC, & OCC data commons infrastructure at the University of Chicago Kenwood Data Center.

# The Gen3 Ecosystem



NHLBI Data Stage

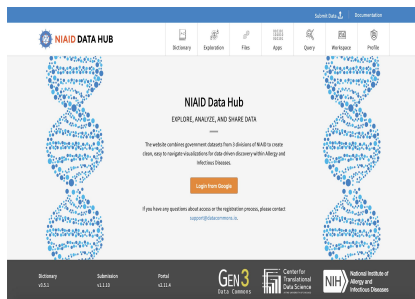NIAID Data Hub

Kids First Data Resource

NHGRI AnVIL

NCI CRDC

Data commons from other foundations

# Narrow Middle Architecture for Data Ecosystem



Diagram: Robert L. Grossman, Progress Towards Cancer Data Ecosystems, The Cancer Journal: The Journal of Principles & Practice of Oncology, 2018, Volume 24, Number 3, May/June 2018.

- Build Notebook in Gen3

- Select virtual cohort in data portal

- Notebook example

- Coming feature for analysis

**GEN3** Data Commons

- Notebooks combine annotation, code, and output visualization



- Gen3 currently supports Jupyter notebooks for a "lightweight workspace"

GEN3 Data Commons

- An authorized user's workspace in a given commons includes a persistent drive in which analysis notebooks, scripts, data files, etc., are saved and persist even after logout
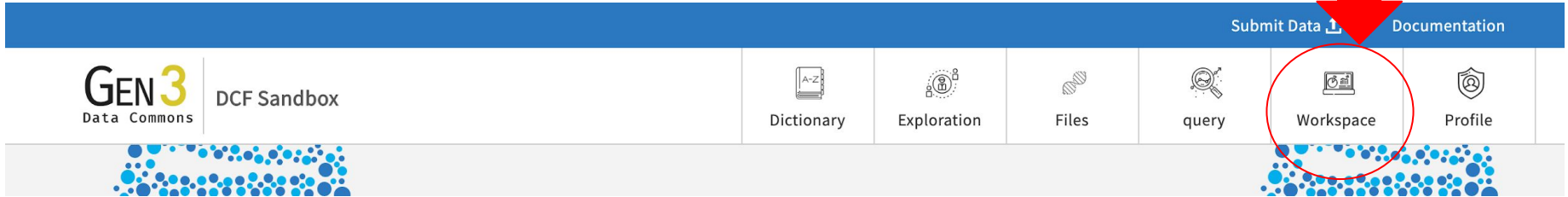
- Gen3 Jupyter notebooks support both R and Python language

- User Documentation available at https://gen3.org/resources/user/analyze-data/

jupyterhub

# Gen3 Python SDK

- The Gen3 SDK facilitates data analysis in notebooks by providing a Python library that makes calls to Gen3 APIs easier.

- Gen3 Python SDK has three classes:
    - Gen3 Auth Helper: Support json web token authentication
    - Gen3 Submission: Submit/Export/Query data from Gen3 submission system
    - Gen3 File Class: Interact with Gen3 file management features

- The Gen3 SDK is available on the python package index (PyPI) at
  https://pypi.org/project/gen3

- For detailed information on how to use the Gen3 SDK, see the Gen3 SDK documentation at http://gen3sdk-python.rtfd.io.

# Access workspace

- Log into commons, select "Workspace"



- Click "Start My Server" to start the Jupyter server in your Workspace

**GEN3** Data Commons

- Choose a virtual machine flavor with the appropriate memory and compute space required for your analysis

- As a Gen3 Data Commons operator, you can configure the different flavors based on the resources available to you, your user community's needs, and what prices you're willing to pay

- The notebook runs a container image that is deployed by kubernetes. The tools and packages in the container are available to anyone selecting the flavor.

## Spawner Options

| ○ | Bioinfo - Python/R 0.5 CPU 256M Mem |
|---|---|
| ○ | Bioinfo - Python/R 1.0 CPU 1024M Mem |
| ○ | Bioinfo - Ariba and Mykrobe 4.0 CPU 15512M Mem |

Spawn

# Docker for notebook

**GEN**³ Data Commons

<> Code    ⓘ Issues **1**    ⑂ Pull requests **2**    ▤ Projects **0**    ▥ Wiki    🛡 Security    �ᴨ Insights

Branch: master ▾    **containers** / **jupyter** / **Dockerfile**    Find file   Copy path

**philloooo** chore(lumpy): add lumpy     0039fbd   26 days ago

**4 contributors**

70 lines (60 sloc)    2.23 KB     Raw   Blame

```
1    # Copyright (c) Jupyter Development Team.
2    # Distributed under the terms of the Modified BSD License.
3    FROM jupyter/scipy-notebook:9e8682c9ea54
4
5    USER root
6
7    RUN pip install --upgrade nbconvert==5.4.1
8
9    # R pre-requisites
10   RUN apt-get update && \
11       apt-get install -y --no-install-recommends \
12       fonts-dejavu \
13       tzdata \
14       gfortran \
15       gcc \
16       libssl1.0.0 \
17       libcurl4-openssl-dev \
18       libssl-dev \
```

```
"jupyterhub": {
    "enabled": "yes",
    "sidecar": "quay.io/cdis/gen3fuse-sidecar:0.1.2",
    "containers": [
        {
            "name": "Bioinfo - Python/R",
            "cpu": 0.5,
            "memory": "256M",
            "image": "quay.io/occ_data/jupyternotebook:1.7.2"
        },
        {
            "name": "Bioinfo - Python/R",
            "cpu": 1.0,
            "memory": "1024M",
            "image": "quay.io/occ_data/jupyternotebook:1.7.2"
        },
        {
            "name": "Bioinfo - Ariba and Mykrobe",
            "cpu": 4,
            "memory": "15512M",
            "image": "quay.io/cdis/niaid-jupyterhub:0.1.1"
        }
    ]
},
```

https://github.com/occ-data/containers

- ## If using an existing notebook and library:

  - Upload any necessary reference files needed for the analysis to your workspace
  - Upload existing Python or R libraries to your workspace
  - You will access clinical data and object files from the data commons within the notebook

# Creating notebook and libraries from scratch in the Workspace

# Prepare your API key for data accessing

- **Create or manage your API keys**
  - API keys are valid for a month
  - Used to receive temporary access token that is valid for only 30 minutes
  - Access token must be sent to Gen3 API to access data in the commons
- **Upload credentials.json to the workspace to allow you to access data within your commons**
  - Be sure your API credentials JSON matches the name of the JSON as you call it in your notebook

# Writing and running Jupyter notebook



Start Writing!

Start Analyzing!

# Tune your source code

- You can stop your notebook to manage your resources responsibly
- If you update your source code or library, you can restart to use the updated code

- Reasons to share a notebook:
  - Review and feedback on methods
  - Other scientists use your analysis on different data based on their access
  - Accompany publication
- Suggestions for sharing notebooks
  - Remove results before sharing
  - Consider GitHub repositories for community accessible notebooks with associated files and libraries
- How to use a shared notebook
  - Upload all libraries and necessary files to your workspace, including notebook
  - Ensure your credentials.json is current and in your workspace

# Virtual cohort selection in data portal

# Virtual cohort selection in data portal

- Select virtual machine flavor with the appropriate memory and compute space required for your analysis

- Import API key and upload or write reference files, library and Jupyter notebook

- Select virtual cohort from data portal and import clinical and object data in virtual machine to run the notebook and tune the library

- Share notebook with community

Outline of the notebook:

- Analyze MRI images getting average cortical thickness measurement in different regions

- Visualize brain surface segmented into different regions

- Compare cortical thickness across groups of patients with different brain disease



Source: https://www.openfmri.org/dataset/ds000030/

Cortical thickness measurement

# Notebook example

Freesurfer Enigma pipeline: recon-all

- Normalize brain signal intensity, skull-stripping, white matter and gray matter segmentation, and delineation of the gray-white interface

- Divide surface into separate cortical regions

- Surface area and mean cortical thickness was extracted for each of the 68 regions (34 per hemisphere)

Source: https://surfer.nmr.mgh.harvard.edu/fswiki/recon-all

```
USER root

COPY license /usr/local/freesurfer/license.txt

ENV FREESURFER_HOME /usr/local/freesurfer
ENV FMRI_ANALYSIS_DIR /usr/local/freesurfer/fsfast
ENV FSFAST_HOME /usr/local/freesurfer/fsfast
ENV FUNCTIONALS_DIR /usr/local/freesurfer/sessions
ENV LOCAL_DIR /usr/local/freesurfer/local
ENV MINC_BIN_DIR /usr/local/freesurfer/mni/bin
ENV MINC_LIB_DIR /usr/local/freesurfer/mni/lib
ENV MNI_DATAPATH /usr/local/freesurfer/mni/data
ENV MNI_DIR /usr/local/freesurfer/mni
ENV MNI_PERL5LIB /usr/local/freesurfer/mni/share/perl5
ENV PERL5LIB /usr/local/freesurfer/mni/share/perl5
ENV SUBJECTS_DIR /usr/local/freesurfer/subjects
ENV PATH $PATH:/usr/local/freesurfer/bin:/usr/local/freesurfer/fsfast/bin:/usr/local/freesurfer/tktools:/usr/local/freesurfer/m

ADD extract_subfields.sh /mnt/
ADD initialize_subDir.sh /mnt/
ADD extract_subcortical.sh /mnt/
ADD outlierDetection.sh /mnt/

RUN apt-get update && apt-get install -y --no-install-recommends curl tar tcsh libglu1-mesa libgomp1 libjpeg62 libxext6 libxtst
  && curl ftp://surfer.nmr.mgh.harvard.edu/pub/dist/freesurfer/6.0.0/freesurfer-Linux-centos6_x86_64-stable-pub-v6.0.0.tar.gz |
  && apt-get install -y --no-install-recommends jq bc libsys-hostname-long-perl && ldconfig && mkdir -p /N/u /N/home /N/dc2 /N/
  && curl "https://surfer.nmr.mgh.harvard.edu/fswiki/MatlabRuntime?action=AttachFile&do=get&target=runtime2012bLinux.tar.gz" -o
  && tar xf /usr/local/freesurfer/runtime2012b.tar.gz -C /usr/local/freesurfer/ \
  && /bin/rm /usr/local/freesurfer/runtime2012b.tar.gz \
  && apt-get remove -y curl \
  && rm -rf /var/lib/apt/lists/*
```

## Spawner Options

○ **Bioinfo - Python/R 0.5 CPU 256M Mem**

○ **Bioinfo - Python/R 1.0 CPU 1024M Mem**
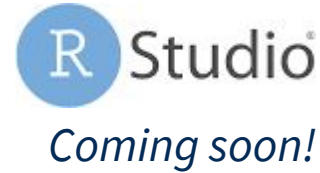
**Brain - Python/R/Freesurfer 1.0 CPU 1.5G Mem**

Spawn

Source: https://surfer.nmr.mgh.harvard.edu

- Now, we will take a look at the Jupyter notebook

GEN3 Data Commons

- Additional tools for the workspace is in development, including R Studio notebooks, Galaxy, and more



*Coming soon!*

- Clinical data export to workspace

- Gen3 workflow execution service. The Gen3 workflow execution service will use its own cwl engine, developed in-house, to execute workflows. User passes the cwl workflow ("packed") as a JSON, as well as a JSON specifying workflow inputs, to the workflow execution service API.

# Learn More

- github.com/uc-cdis

- gen3.org

- Gen3 Community on Slack

- support@datacommons.io

- ctds.uchicago.edu

Selected Data Commons Using Gen3

Questions?