# A clinical data mesh for quality improvement and research in healthcare

**Mike Hogarth**
MD, Clinical Research Information Officer, UC San Diego Health

**Tom Covington**
CEO of Tag.bio

**Mark Mooney**
VP Customer Success, Tag.bio

Campus
L>SA

# Agenda (50 minutes)

| | |
|---|---|
| **UCSD and Tag.bio history** | 15 min |
| **What is a Data Mesh?** | 10 min |
| **Nightingale demo** | 20 min |
| **Q & A** | 5 min |

Campus L>SA

# UCSD

**Mike Hogarth,** MD, Clinical Research Information Officer, UC San Diego Health

Campus L>SA

# The problem (use case) we are trying to solve

- Provide clinical data (protected health information - PHI) to UCSD biomedical researchers:
  - Securely
  - Swiftly
  - Simply
  - Standardized
  - Semantically consistent

Campus LSA

# Data (and analyses) should be FAIR.



| **F** | **A** | **I** | **R** |
|---|---|---|---|
| Findable | Accessible | Interoperable | Reusable |
| "I know where all our data is" | "I can access any of the data that I need" | "I use one language for all my requests" | "I use the existing data to answer new questions" |

Campus
L>SA

# UCSD Health Secure Research Cloud in AWS iDASH

- What is it?
  - Integrating **D**ata for **A**nalysis, anonymization, and **SH**aring (iDASH) - v2.0
  - A secure computing environment for sensitive data
  - HIPAA compliant
  - Health system CSO approved for use involving protected health information (PHI)
  - Includes virtual research desktops and VMs in a locked-down AWS VPC without the ability for users to connect to the "outside"

- Why we needed it
  - c2017 - EHR data (PHI) for research routinely given to investigators through a download -- some data ended up shared it with external entities without data use agreements, no visibility into the location/use of the data, the ADCS situation along with ~800 other AWS accounts by UCSD Health staff/faculty without visibility or controls

Campus
L>SA

# UCSD health secure research cloud with nodes in a VPC

## THE VIRTUAL RESEARCH DESKTOP (VRD)

- It is a modified version of the Amazon Web Services (AWS) Windows 10 "Workspace" virtual machine
- Runs in the protected UCSDH Secure Cloud in AWS
  - in the AWS HIPAA environment
  - approved by UCSDH CSO for PHI
- Provisioned with:
  - SPSS
  - R/RStudio
  - Python/PyCharm
  - Java 8 JDK
  - Depending on approval, access to internal databases – ie, UC CORDS
  - tag.bio based access to available databases



UCSD Health Virtual Research Desktop

**UC San Diego**
ALTMAN CLINICAL AND TRANSLATIONAL
RESEARCH INSTITUTE

NIGHTINGALE

DIAGRAM of the CAUSES of MORTALITY
IN THE ARMY IN THE EAST.

2.
APRIL 1855 to MARCH 1856.

1.
APRIL 1854 to MARCH 1855.

The Areas of the blue, red, & black wedges are each measured from
    the centre as the common vertex.
The blue wedges measured from the centre of the circle represent area
    for area the deaths from Preventible or Mitigable Zymotic diseases, the
    red wedges measured from the centre the deaths from wounds, & the
    black wedges measured from the centre the deaths from all other causes.
The black line across the red triangle in Nov.r 1854 marks the boundary
    of the deaths from all other causes during the month.
In October 1854, & April 1855, the black area coincides with the red;
    in January & February 1856, the blue coincides with the black.
The entire areas may be compared by following the blue, the red & the
    black lines enclosing them.

**Florence Nightingale**
OM RRC DStJ



Florence Nightingale, c. 1860

Campus L>SA

# History

Tom Covington, CEO Tag.bio

# How did the data mesh platform arise?

- Jesse Paquette (CSO Tag.bio) worked at UCSF Helen Diller, Family Comprehensive Cancer Center (2007-2010)

    - Working with Oncology researchers he realized that enabling them to answer their own questions would speed the turnaround time of question to answer

    - Created an initial software called EGAN (Exploratory Gene Association Networks)

    - Formed Tag.bio in 2014 with Tom Covington (CEO) and built first versions of what were then called Flux Capacitors or FC's but became data nodes.

    - Began projects with UCSF Med Center on billing, encounters and claims data in 2018.

    - Realized that the architecture we were working on was an implementation of a "data mesh" after reading Zhamak Dhegani's article: How to Move Beyond a Monolithic Data Lake to a Distributed Data Mesh

# What's The History?

- **In January 2020 met Tag.bio at Precision Medicine World Conference**

  

  - Initial discussion about work with UCSF on value based healthcare.

  - Set up a visit to UCSD in February

  - Initiated research collaboration with Mike Hogarth in March

- **The Pandemic**

  

  - Realized there was an immediate need to make COVID data accessible

  - Built first COVID registries in April



**Building a range of patient registries at the present time**

# What is a Data mesh?

Mark Mooney, VP of Customer Tag.bio

Campus
L>SA

# Data nodes deployed and registered in a decentralized
# Data Mesh



Similar to an app store
or a library of data
products.

Zhamak Dehghani (Thoughtworks): How to Move Beyond a Monolithic Data Lake to a Distributed Data Mesh

# UCSD health secure research cloud with nodes in a VPC

# Tying 3 technical components...



## Data

## Algorithms

## Visualization

...into a domain-driven data node

**Smart API**

**Data Map**

**Algorithms**

**Data node**

Campus LYSA

# NODE: Data Map

## Types

- Proprietary data (e.g. patients, billing, etc.)
- Public data (e.g. TCGA, CDC, etc.)
- Data services (e.g. annotation, usage, etc.)
- **Emerging data types**

## Sources

- Siloed data
- Data lake
- Data node

## Formats

- CSV
- JSON
- SQL

Campus LYSA



# UC CORDS COVID-19



Combining healthcare data from across the six University of California medical schools and health systems

Aug 14 2020
- 175,517 COVID tested patients
- 6,056 COVID+ patients
- all labs, meds, vitals, 29 ICU data elements
- 319,952,837 "data points"

**de-identified data
Ingested in 5 hours**

# NODE: Algorithms



**Classic statistical methods**

PCA, tSNE, UMAP,
t-test, Mann-Whitney,
Hypergeometric,
Cox regression,
Paired analysis over time,
Pathway analysis, etc.

**Integration with**

R

python

Machine Learning / AI

Campus LYSA

# NODE: Smart API/UX

**Analysis Apps**
**Made** by the data scientist
**Used** by the researcher

Smart API

Data Map

Algorithms

N-Month Survival

Cox Survival

Expression

Data Quality

PCA

Expression Cutoff

DNA Mutation/Variation

Data node

Data node

Data node

Campus LYSA

NODE: Fast MVP robust CI/CD

NODE: DIY

Project | python_elastic_net_cross_validation.json

- fc-topaz
  - .git
  - .idea
  - config
  - docker-scripts
  - ipy
  - protocols
    - annotations
    - argument_references
    - argument_sets
    - arguments
    - callbacks
    - cohort_protocols
    - entity_annotation
    - variable_protocols
    - variables
    - batch_rich_elastic_net_multiple_regression.json
    - download.json
    - elastic_net_multiple_regression.json
    - multiple_regression.json
    - numeric_focus.json
    - python_create_correlations_download.json
    - python_create_correlations_embed.json
    - python_elastic_net_cross_validation.json
    - rich_elastic_net_multiple_regression.json
    - variable_summary.json
  - publications
  - python
  - scripts
  - shell_scripts
  - .gitignore
  - .rancher-pipeline.yml
  - deployments.yml
  - protocols.txt
  - README.md
  - scripts.txt

```json
{
  "protocol_definition": {
    "name": "elastic_net_crossvalidation",
    "visible": true,
    "thumbnail": "https://tag-client-images.s3-us-west-2.amazonaws.com/protocol-thumbnails/fc-topaz/python.png",
    "title": "Elastic net cross validation",
    "description": "This protocol performs elastic net cross validation.",
    "category": "Python plugins",
    "argument_sets": [
      "protocols/argument_sets/background_cohort_argument_set.json",
      "protocols/argument_sets/model_outcome_argument_set.json",
      "protocols/argument_sets/model_cross_validation_input_argument_set.json",
      "protocols/argument_sets/elastic_net_cross_validation_parameters_argument_set.json"
    ],
    "tags": [
      "Python",
      "Elastic net",
      "Cross validation"
    ]
  },
  "method": "external",
  "sdk": "python",
  "plugin": "python/elastic_net_cross_validation.py",
  "output_type": "html",
  "background": "protocols/argument_references/background_cohort_reference.json",
  "analysis_variables": [
    "protocols/variables/auto/patient_id_collection.json",
    "protocols/argument_references/model_outcome_variable_reference.json",
    "protocols/argument_references/model_input_variables_reference.json",
    "protocols/argument_references/plos_one_model_cv_input_variables_reference.json"
  ]
}
```

protocols/python_elastic_net_cross_validation.json    18:25       LF    UTF-8    JSON    ⑂ master    ⟳ Fetch    ⬡ GitHub    Git (1)

# Nightingale registries
## (a de-identified, automated, OMOP data product)

Mark Mooney, VP of Customer Tag.bio

UC CORDS Research Registry - All Tested Patients History View

**175517 Patients**

University of California CORDS registry of all patients, both positive and negative for COVID-19.

Server version: 2.34.18
Aug 26, 2020
9:02:06 AM

UC BRAID CORDS

Type here to filter datasets

☐ CORDS - UC COVID Patient Registry (2)
☐ UCSD COVID Patient Registry (2)
☐ Services (3)
☐ Synthetic COVID-19 Patient Data (1)
☐ The Cancer Genome Atlas (TCGA) Datasets (4)

## UC CORDS Research Registry - Positive Patient History View

**6056 Patients**

University of California CORDS registry of COVID-19 positive patients.

Server version: 2.34.18
Aug 26, 2020
9:02:08 AM

UC BRAID CORDS

### Synthetic COVID-19 Patient Data

## Synthetic COVID-19 demo data

**1835 patients**

Using synthetic, OMOP mapped, COVID 19 data to enable analysis demonstration and review.

Server version: 2.34.18
Aug 26, 2020
12:56:18 PM

### The Cancer Genome Atlas (TCGA) Datasets

## TCGA Pan-Cancer Atlas and the Immune Landscape of Cancer

NIGHTINGALE

## COVID-19 Comparison Apps

| Comparison of COVID-19 Positive Patients | Comparison of Specific Variables for COVID-19 Positive Patients | Comparison of COVID-19 Positive Patients with Pre-Existing Conditions |
|---|---|---|

Type here to filter protocols

☐ Comparison
☐ Download
☐ Overview
☐ Specialty
☐ Summary

## COVID-19 Specialty Apps

**Cox Survival**

**Cox Survival Using Specific Variables**

This protocol allows you to perform cox survival analysis using specific variables on a COVID-19 positive patient cohort that you define.

**Click to configure**

## Data Download Apps

**Patients Data Download**

Home / Select a dataset / Synthetic COVID-19 demo data

## Synthetic COVID-19 demo data

**1835 patients**

Using synthetic, OMOP mapped, COVID 19 data to enable analysis demonstration and review.
View reference (new tab)

Protocols          Subjects          Run a script

**Overview Apps**

Type here to filter protocols

Overview of Data

☐ Comparison
☐ Download
☐ Overview
☐ Specialty
☐ Summary

**COVID-19 Summary Apps**

Summary of COVID-19 Positive Patients

Summary of Specific Variables for COVID-19 Positive Patients

Summary of COVID-19 Positive Patients with Pre-Existing Conditions

**COVID-19 Comparison Apps**

**Comparison of COVID-19 Positive Patients**

This protocol compares clinical outcomes for all COVID-19 positive patients with a defined sub-cohort of COVID-19 positive patients.

**Click to configure**

**Comparison of Specific Variables for COVID-19 Positive Patients**

**Comparison of COVID-19 Positive Patients with Pre-Existing Conditions**

**COVID-19 Specialty Apps**

**Cox Survival**

**Cox Survival Using Specific Variables**

Type here to filter protocols

☐ Comparison
☐ Download
☐ Overview
☐ Specialty
☐ Summary

**Cohort builder**

🏥 Visit Details                                                                        ˅

👤 **Patient Details**                                                                  ˄

Clear all

Gender                                                                    FEMALE  ˄
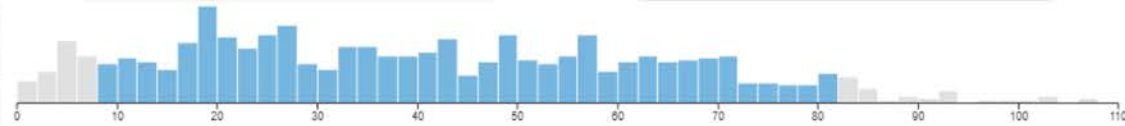☐ Select all
☑ FEMALE
☐ MALE

Ethnicity                                                                              ˅

Race                                                                                   ˅

Age range                                                                              ˄

40                                   ✕              80                             ✕



*412 patients*     Cancel     **Save & use**

**Cohort builder**

Please specify a name for this cohort:

f-40 to 80

Selected parameters

**Patient Details**
Gender: FEMALE
Minimum age: 40
Maximum age: 80

412 patients

Back     Cancel     Use now     Save

# NIGHTINGALE

Home / Select a dataset / Synthetic COVID-19 demo data / Comparison of COVID-19 Positive Patients / Configure and run protocol

## Comparison of COVID-19 Positive Patients

This protocol compares clinical outcomes for all COVID-19 positive patients with a defined sub-cohort of COVID-19 positive patients.
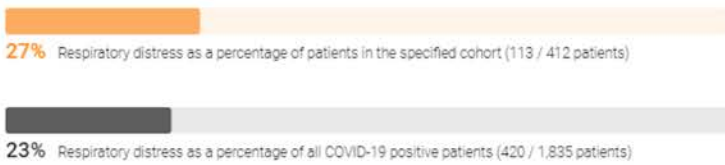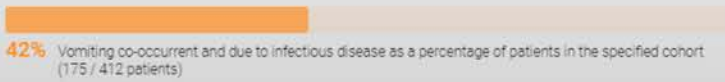
**› Instructions**

COVID-19 Patient Details

### 🌸 COVID-19 Positive Cohort   ⌃
*Select additional criteria to further define the COVID-19 positive cohort.*

Clear all

COVID-19 Positive cohort                                    f-40 to 80

🔧 **Rebuild cohort**     ☣ **Use saved cohort**

Collections of Variables to Analyze

### ✳ Data from Entire Patient Record to Analyze   ⌃
*Select data to analyze.*

Clear all

Outcomes recorded only after COVID positive status     **Multiple values** ⌃

☐ Select all          Filter values...

☑ Condition - After COVID Positive Diagnosis
☐ Drug - After COVID Positive Diagnosis
☑ Measurement - After COVID Positive Diagnosis
☑ Observation - After COVID Positive Diagnosis
☐ Procedure - After COVID Positive Diagnosis

### Selected parameters   ⌃

**COVID-19 Positive Cohort**
COVID-19 Positive cohort:   f-40 to 80

**Data from Entire Patient Record to Analyze**
Outcomes recorded only after COVID positive status:
Condition - After COVID Positive Diagnosis,
Observation - After COVID Positive Diagnosis,
Measurement - After COVID Positive Diagnosis

**Sort Results**
Sort attribute:   Tag.score
Sort direction (default >):   >
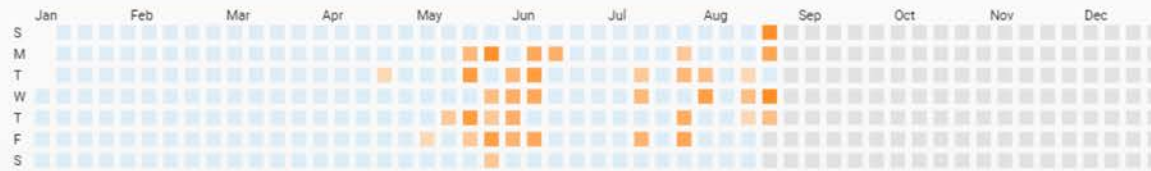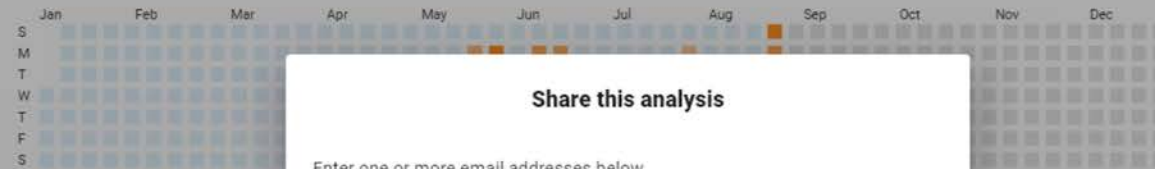Define a maximum number of results to view across
all analyzed variables.:   2500

▶ **Run protocol**

**Cohort builder**

🔵 Demographics                                                          ⌄

🟡 Cancer                                                                ⌄

🔴 Tumor event                                                           ⌄

🔵 Patient and sample IDs                                                ⌄

🟢 Survival                                                              ⌄

🔴 Tissue                                                                ⌄

🟠 Histology                                                             ⌄

🟣 Stage                                                                 ⌄

🔵 Immune subtypes                                                       ⌄

🔵 Copy number variants                                                  ⌄

*10967 samples*    Cancel    **Save & use**

# SCGB2A2

☰ Figure          ⁖ Drill down          🗐 Annotation

tag.score          Expression
☐    **750**        **SCGB2A2**                                    ➡  🔗  📊
       p = 0        `Higher than expected`

**9.51**  Average expression of SCGB2A2 for patients in the focus cohort (1082 patients)

**1.52**  Average expression of SCGB2A2 for all patients in the background cohort (8414 patients)

**0.342**  Average expression of SCGB2A2 for patients NOT in the focus cohort (7332 patients)

Dimension 1

trp.ucsd.edu/fc-pan-can/umap/results

# SCGB2A2

≡ Figure     ⠿ Drill down     📋 Annotation

## Annotation

**Chromosome** : 11
**Cytoband** : 11q12.3
**Ensembl RNA ID** : ENST00000227918.3, ENST00000525380.1
**Ensembl gene ID** : ENSG00000110484
**Ensembl protein ID** : ENSP00000227918.2, ENSP00000431997.1
**GO PubMed** : 25416956, 21873635
**GO biological process** : GO:0030521: androgen receptor signaling pathway, GO:0008150: biological_process
**GO cellular component** : GO:0005615: extracellular space, GO:0005575: cellular_component
**GO molecular function** : GO:0003674: molecular_function, GO:0005515: protein binding
**Gene ID** : 4250
**Gene name synonyms** : mammaglobin A
**Gene symbol** : SCGB2A2
**Gene synonym** : UGB2, MGB1, PSBP1
**Gene type** : protein-coding
**GeneBank accession** : AF015224
**HGNC previous symbol** : MGB1, PSBP1
**HGNC symbol** : SCGB2A2
**HGNC synonym** : MGC71974, UGB2
**Locus group** : protein-coding gene
**Locus type** : gene with protein product
**Modification date** : 2019/10/12
**Mouse genome database ID** : MGI:3780828
**Orientation** : Positive strand
**Other database ID** : HGNC:HGNC:7050, MIM:605562, Ensembl:ENSG00000110484
**Other gene designation** : prostatic steroid binding protein 1, mammaglobin-A, mammaglobin 1
**PubMed ID** : 17192791, 18251583, 27477018, 25416956, 21411781, 21976532, 15489334, 22994369, 17653857, 22897908, 26276775, 16110760, 22963676, 24823311, 20586026, 17071045, 21744998, 18846421, 16203799, 20092039, ... (44 more)
**Refseq RNA** : XM_005274005.3, NM_002411.4
**Refseq gene** : SCGB2A2

# Differential expression

Variable collection
**Expression**

Type here to filter 100 results

| tag.score | Expression | | |
|---|---|---|---|
| **750**  p = 0 | **SCGB2A2**  Higher than expected | ➡ | 🔗 |

**9.51**  Average expression of SCGB2A2 for

**1.52**  Average expression of SCGB2A2 for

**0.342**  Average expression of SCGB2A2 for patients NOT in the focus cohort (7332 patients)

| tag.score | Expression | | |
|---|---|---|---|
| **750**  p = 0 | **PIP**  Higher than expected | ➡ | 🔗 |

## Download as csv

tagbio-analysis-export.csv

**Download**    Cancel

Dimension 1

NIGHTINGALE

1835 patients

Using synthetic, OMOP mapped, COVID 19 data to enable analysis demonstration and review.

Server version: 2.34.18
Aug 26, 2020
12:56:18 PM

Type here to filter datasets

☐ CORDS - UC COVID Patient Registry (2)
☐ UCSD COVID Patient Registry (2)
☐ Services (3)
☐ Synthetic COVID-19 Patient Data (1)
☐ The Cancer Genome Atlas (TCGA) Datasets (4)

## The Cancer Genome Atlas (TCGA) Datasets

### TCGA Pan-Cancer Atlas and the Immune Landscape of Cancer

10967 samples

Combined data from 33 cancer types from the 2018 TCGA Pan Cancer Clinical Data Resource

Server version: 2.34.18
Aug 27, 2020
7:03:24 AM

### METABRIC Breast Cancer

1980 patients

A dataset with clinical and multi-omics data for 1980 breast cancer patients (METABRIC, Nature 2012 & Nat Commun 2016).

Server version: 2.34.18
Aug 26, 2020
7:46:57 PM

### Head and Neck Cancer (TCGA)

530 samples

A dataset with clinical and multi-omics data for 530 head and neck cancer patients (from TCGA).

Server version: 2.34.18

Tag.bio - Select a dataset

trp.ucsd.edu/fc

NIGHTINGALE

**UC CORDS Research Registry - Positive Patient History View**

6056 Patients

University of California CORDS registry of COVID-19 positive patients.

Server version: 2.34.18
Aug 26, 2020
9:02:08 AM

UC BRAID CORDS

Type here to filter datasets

☐ CORDS - UC COVID Patient Registry (2)
☐ UCSD COVID Patient Registry (2)
☐ Services (3)
☐ Synthetic COVID-19 Patient Data (1)
☐ The Cancer Genome Atlas (TCGA) Datasets (4)

**Synthetic COVID-19 Patient Data**

**Synthetic COVID-19 demo data**

1835 patients

Using synthetic, OMOP mapped, COVID 19 data to enable analysis demonstration and review.

Server version: 2.34.18
Aug 26, 2020
12:56:18 PM

**The Cancer Genome Atlas (TCGA) Datasets**

**TCGA Pan-Cancer Atlas and the Immune Landscape of Cancer**

10967 samples

Combined data from 33 cancer types from the 2018 TCGA Pan Cancer Clinical Data Resource

Server version: 2.34.18
Aug 27, 2020
7:03:24 AM

# Summary

Mike Hogarth MD, Clinical Research Information Officer, UC San Diego Health

Campus L>SA

# NIGHTINGALE

## A tool for data exploration and analysis

- We have installed the tag.bio system in our research cloud and it has access to data sets in our 'secure data commons database'

- The Nightingale portal provides population level access and ability to perform analysis

- A user can 'slice' the cohort and select specific analyses (demographic, survival, comparison between cohorts)

- Planned, pending approval, provide 'download' of limited data set (LDS) row-level data from selected data set into the investigator's virtual research desktop for further analysis

Campus LYSA

# Help us evolve the mesh

- What other registries should be available?
- How would you like to query them?
- Are there public data sources you would like to see here?
- Could we use the mesh for other data sources?

Please contact Mike Hogarth at mihogarth@health.ucsd.edu with suggestions or comments.

Campus L>SA

# Thank You!

# Questions?

# Next presentation

**Come see our next UC TECH Presentation 9/03:**

Email Overload: Practical Tools for Influencing Email Volume in the Age of Telecommuting

Are you overwhelmed by the number of emails you receive daily? Has email management become a burdensome core task that monopolizes your time? When volume exceeds 200+ emails a day, generic email tips & tricks for email management simply won't cut it. This session will go beyond email platform use to focus on email management from a behavior modification and process improvement perspective. We will cover practical tools and strategies for actively managing virtual work and interactions with your co-workers more effectively, giving you the ability to actually influence the volume of emails you receive.

Speaker:

Loralyn Cross, Office of Research Affairs, UC San Diego

Campus L>SA

# Reference slides

# Enabling doctors to provide instant answers

## 9 years of billing and encounter data

- All inpatient and outpatient data in a combined dataset

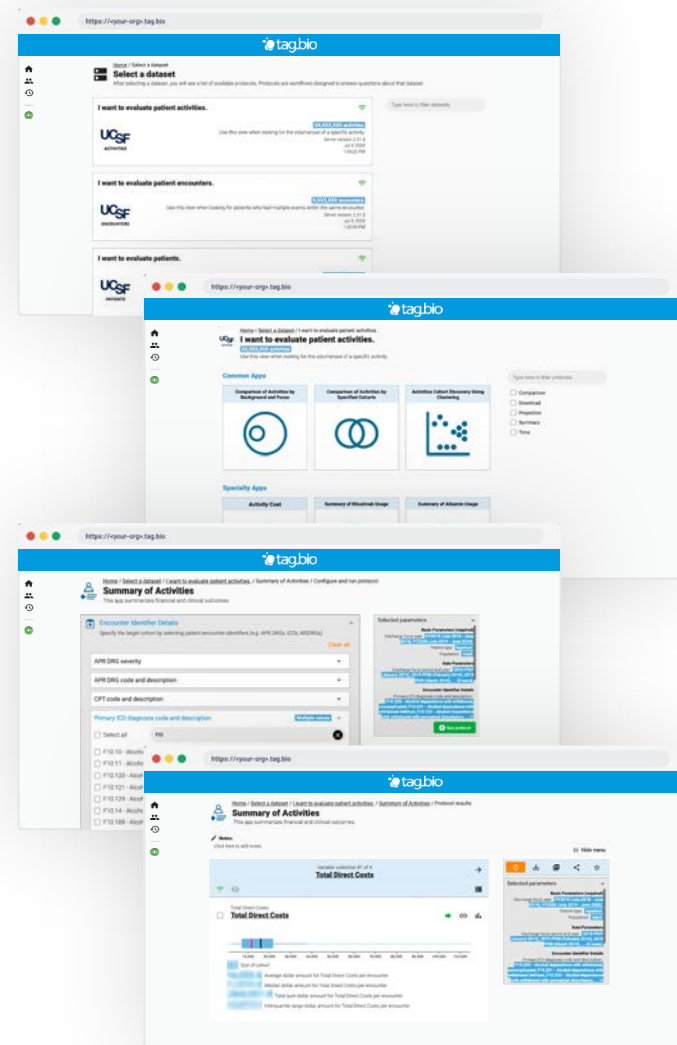➔ Over **2,000 analyses** performed by value improvement physicians in the past year

"The ability to have this kind of **on-demand information** completely **changes the culture**.
I can't imagine doing my job without the Tag.bio platform."

*- Jahan Fahimi, Director of Value Improvement at UCSF Health*

**Jahan Fahimi, MD, PhD,**
Associate Professor of Emergency Medicine,
Director of Value Improvement at UCSF Health

Campus LYSA

# What does this enable?

|  | **Data** | **Development** | **Collaboration** |
|---|---|---|---|
| **Node** | Is a Data Product | Build a rapid Data Product | SMART API |
|  | Domain-specific analysis | Iterative dev cycle | Secure data |
|  | Immutable data and transient node | Integrate other functions (R, Python, ML) | Secure deploy |
|  | Deploy anywhere | | |
| **Mesh** | Many Nodes, many data types | Transfer apps between nodes | Distributed querying |
|  | Node functional diversity | Federated functionality | Centralized analysis |
|  | Distributed analyses | No down time | Public/Private nodes |
|  | Network effect on data value | An ecosystem of nodes | |
| **Portal** | Publish data with analyses | Rapidly populated by an admin | Reproducible, replayable analyses |
|  | Track analyses through history | Cohorts into Nodes | Share analyses |
|  | Create COHORT/UDAT | Transactional Nodes | Versioned resources |
|  | Reference & Annotation nodes | Usage allows evolution of mesh | De-silo analysis types |

LYSA