# 9PM

2AM in London (GMT), 11AM in Tokyo (GMT+9)

## Multiscale Models

**Moderator:** Katy Börner, *Indiana University*
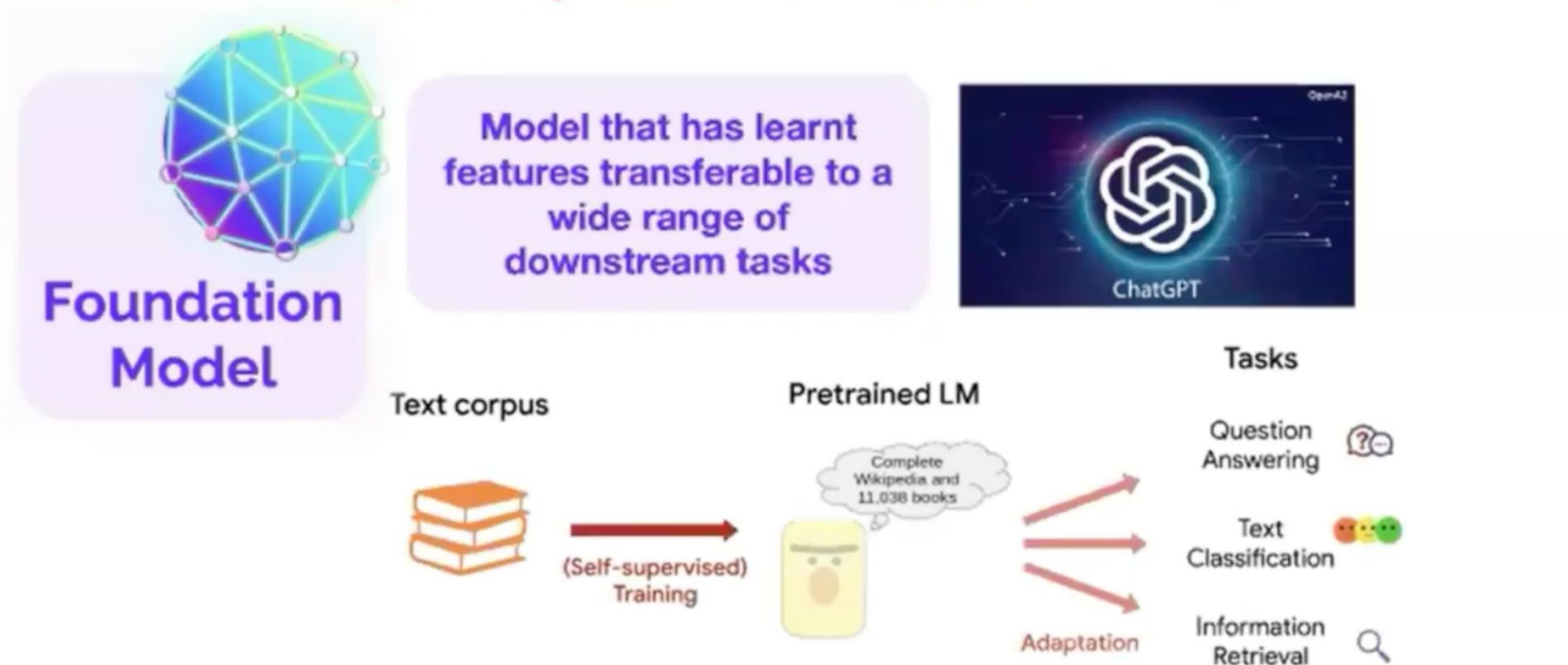
**Presenters:**

- Maria Brbic, *Swiss Federal Institute of Technology Lausanne, Switzerland*
- Filipi N. Silva, *Indiana University*

# Maria Brbic, *Swiss Federal Institute of Technology Lausanne*

# AI Revolution

## Generative AI paradigm and the era of foundation models



**Foundation Model**

Model that has learnt features transferable to a wide range of downstream tasks

ChatGPT

Text corpus → (Self-supervised) Training → Pretrained LM (Complete Wikipedia and 11,038 books) → Adaptation → Tasks

Tasks:
- Question Answering
- Text Classification
- Information Retrieval

How can we leverage these AI advances in single cell biology?

What are their current limitations for biomedical applications?

# Single-cell Data Is Challenging for Today's AI

**1** Heterogenous experiments

**2** Novel and unknown phenomena

**3** Different modalities with different challenges

1 Heterogenous experiments
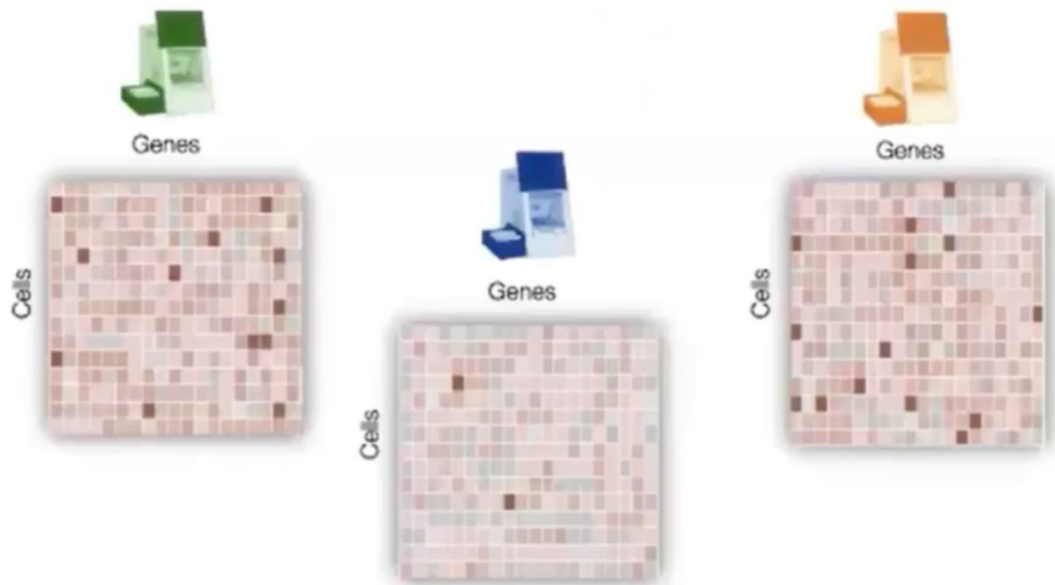
2 Novel and unknown phenomena

3 Different modalities

Today's talk: How to overcome some of these challenges

# On Heterogeneity

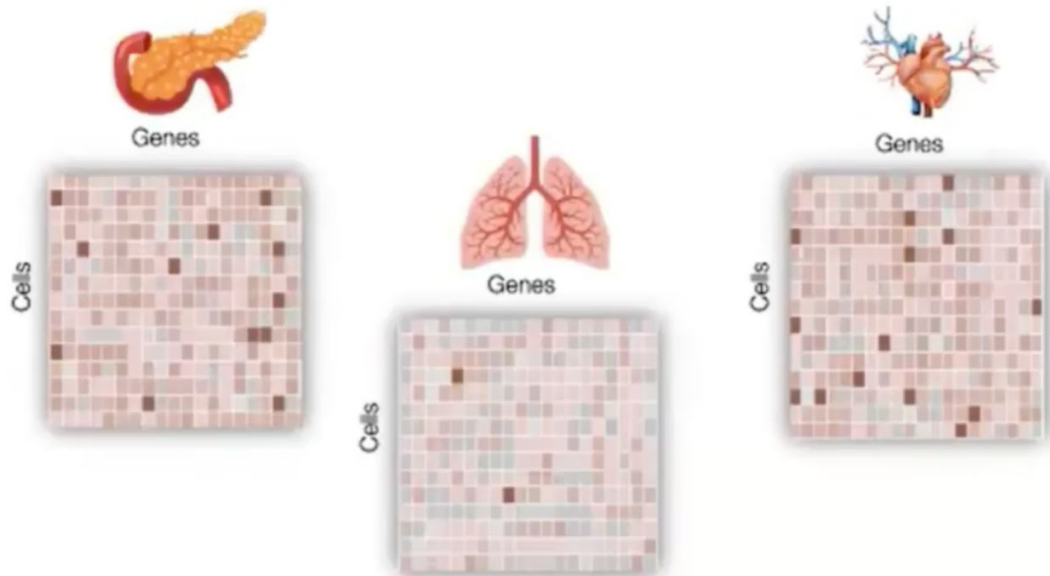## Discovering Cell Types Across Tissues, Disease States & Species

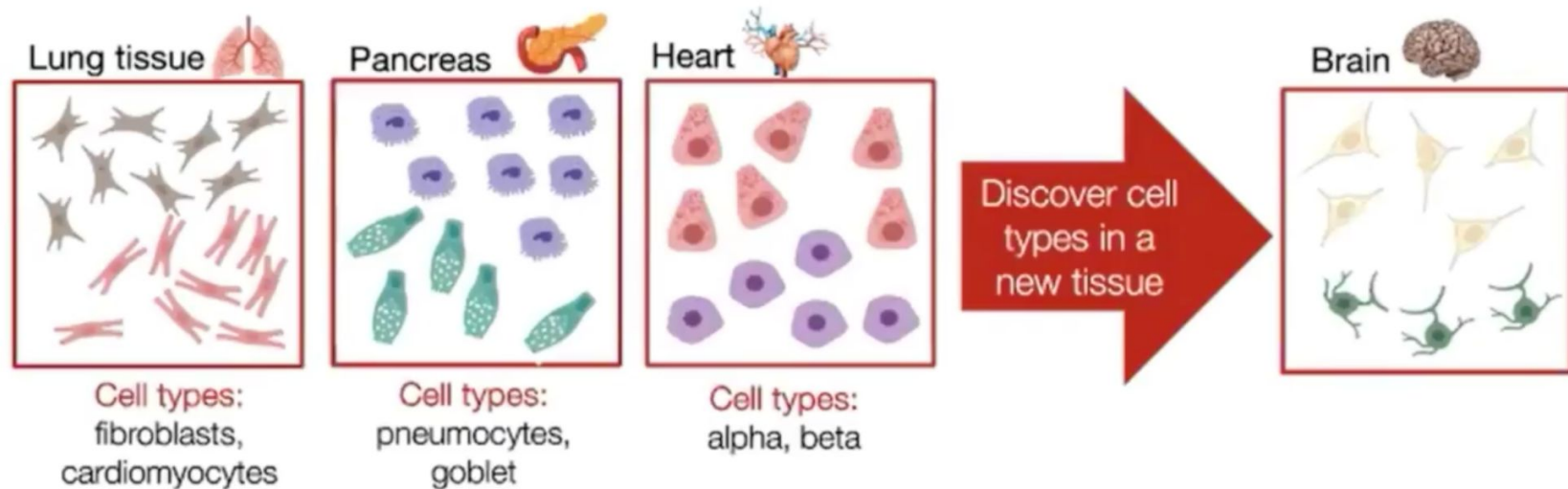# Data with Large Heterogeneity

## different labs…

# Data with Large Heterogeneity

## different tissues…

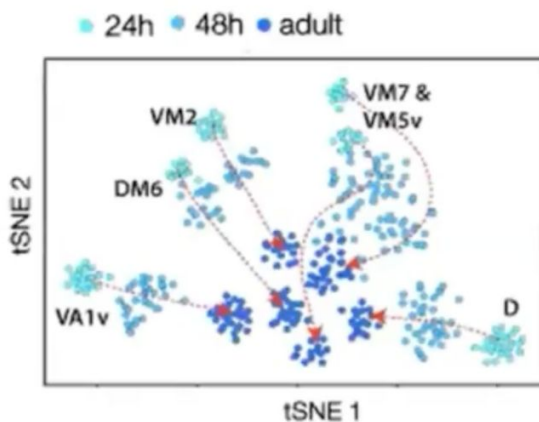How do we jointly analyze and gain new insights from these heterogenous datasets?

# MARS: Learn Cell Embeddings to Discover Novel Cell Types

Lung tissue

Cell types:
fibroblasts,
cardiomyocytes

Pancreas

Cell types:
pneumocytes,
goblet

Heart

Cell types:
alpha, beta

Discover cell types in a new tissue

Brain

Brbic et al. *Nature Methods* '20

# Cell Type Discovery across Experiments



Across tissues of the
Mouse Cell Atlas
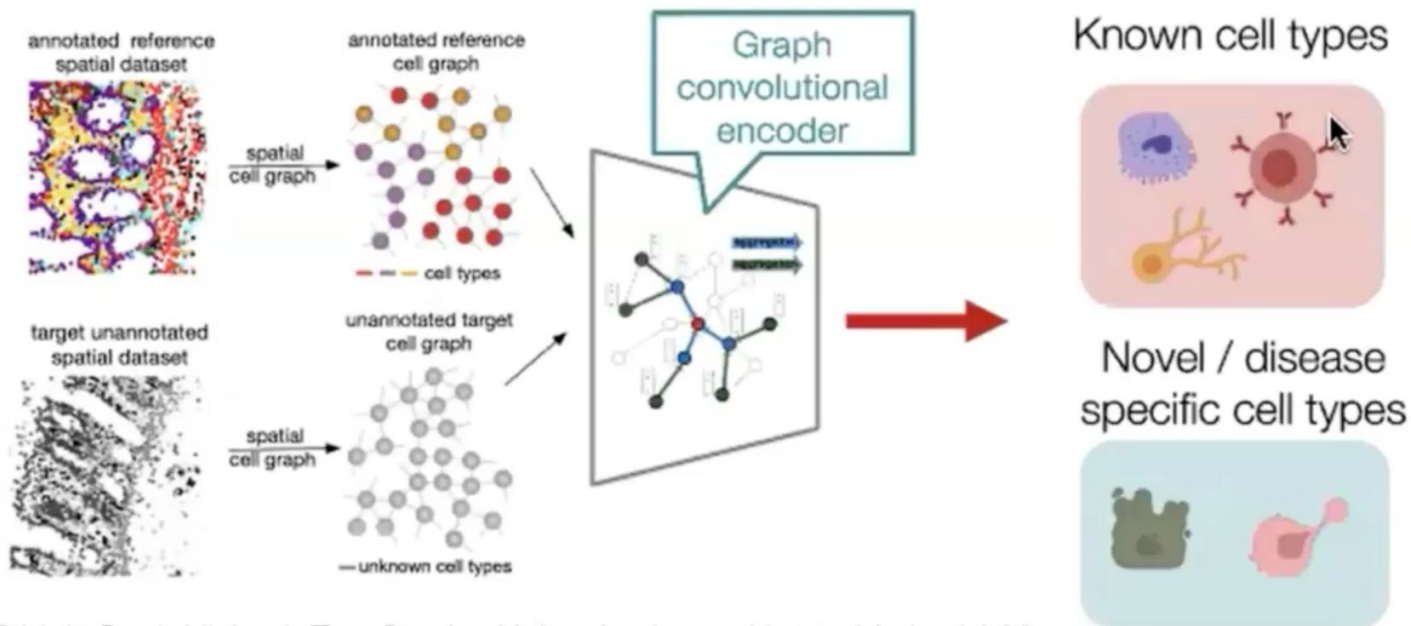
Xie*, Brbic* et al. *eLife* '21

Fly Cell Atlas

Li*, Janssens* et al. *Science* '22

14

STELLAR: Novel Cell Type Discovery Across Conditions

Brbic*, Cao*, Hickey*, Tan, Snyder, Nolan, Leskovec *Nature Methods* '22

# Towards Universal Cell Embeddings

Can we create cell embeddings for any species, any set of genes?

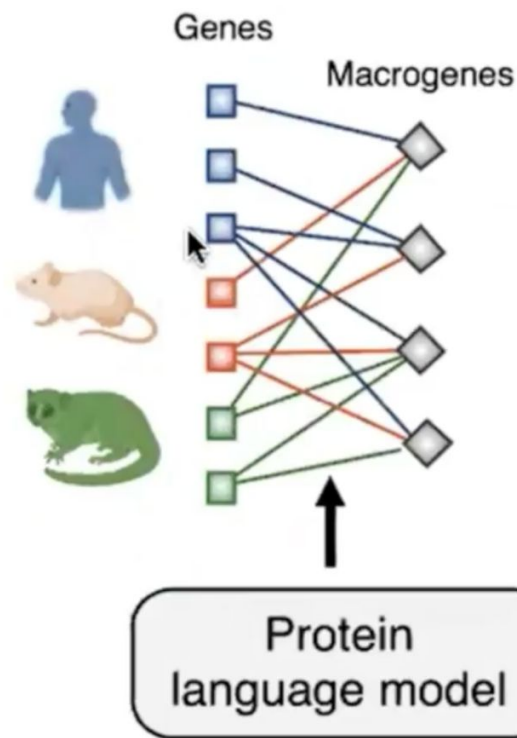Tabula Muris
Nature '18, '20

Fly Cell Atlas
*Science '22 '23*

Tabula Sapiens
*Science '22*

# SATURN: Integrating Datasets across Species
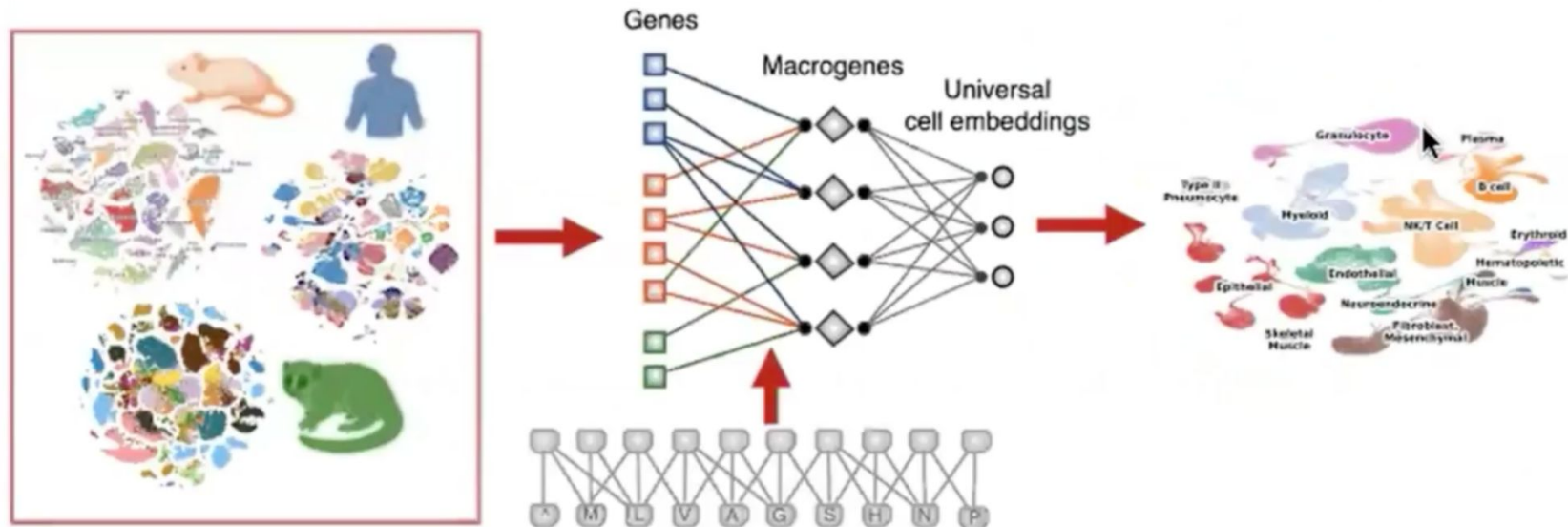
Genes
Macrogenes
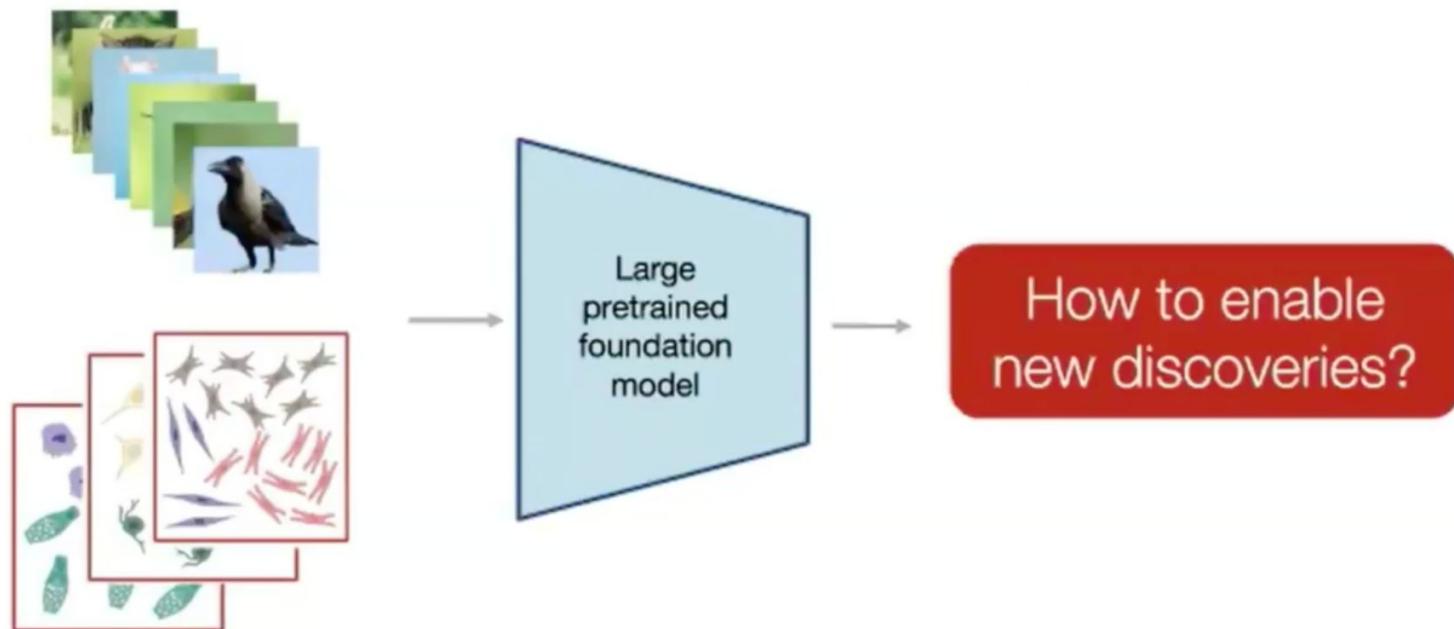Universal cell embeddings

Rosen*, Brbic*, Roohani*, Swanson*, Li, Leskovec *Nature Methods'* 24

# On Discovery

Enabling Discovery from Foundations Models

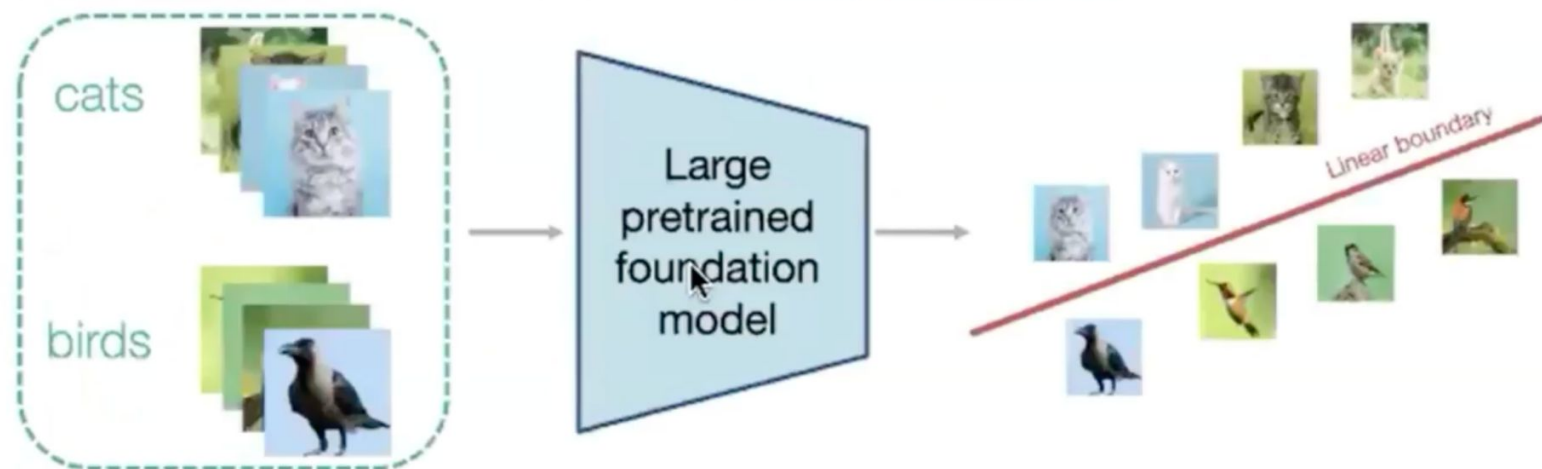# How To to Enable New Discoveries from Foundation Models?



Large pretrained foundation model

How to enable new discoveries?

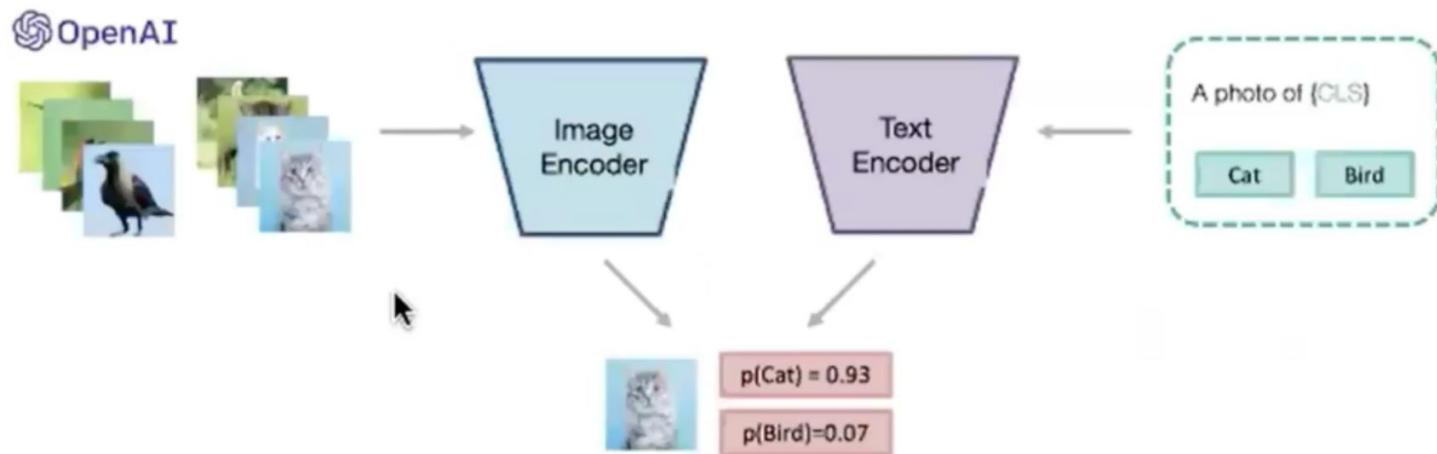# Current Paradigms Still Require Supervision

- Current paradigms:

  1. Fine-tune on the task of interest using labeled data



Maria Brbić, EPFL
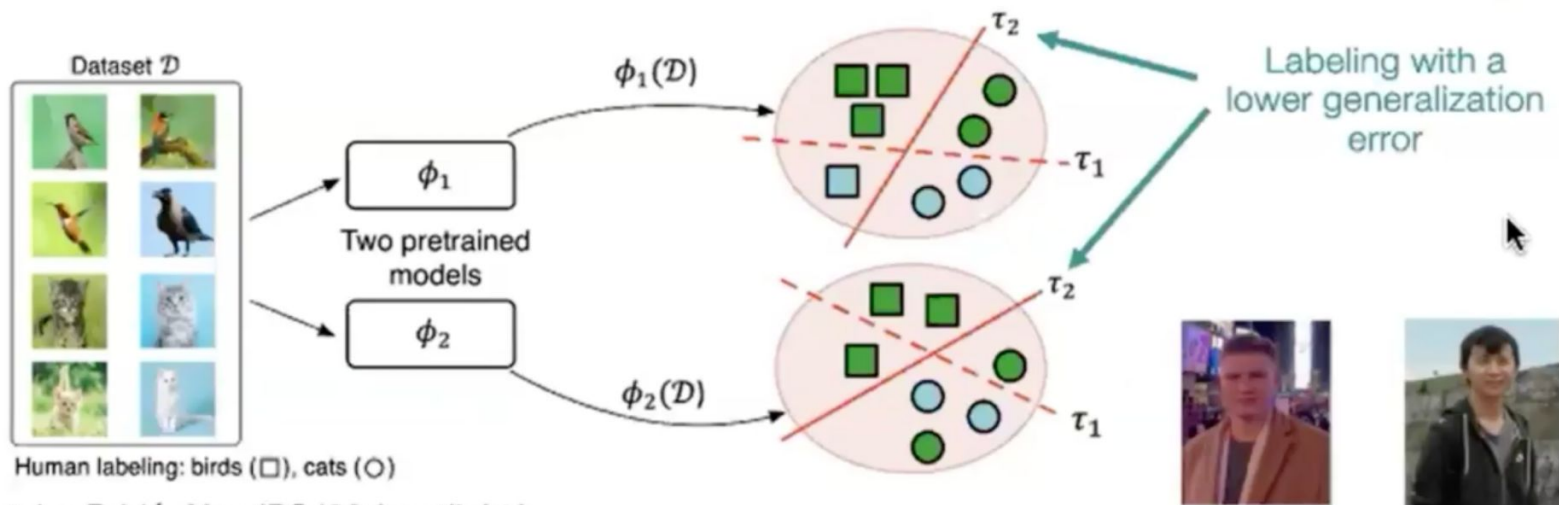
# Current Paradigms Still Require Supervision

- Current paradigms:
  2. Zero-shot transfer on the task of interest using instructions

# How to Infer Labeling without Any Supervision?

**Key idea:** Search for a labeling such that linear models will generalize well in different representation spaces



Dataset $\mathcal{D}$

$\phi_1(\mathcal{D})$

Two pretrained models

$\phi_1$

$\phi_2$

$\phi_2(\mathcal{D})$

Human labeling: birds (□), cats (○)

$\tau_2$

$\tau_1$

Labeling with a lower generalization error

Gadetsky, Brbić. *NeurIPS* '23 (spotlight )
Gadetsky*, Jiang*, Brbić. *ICML* '24

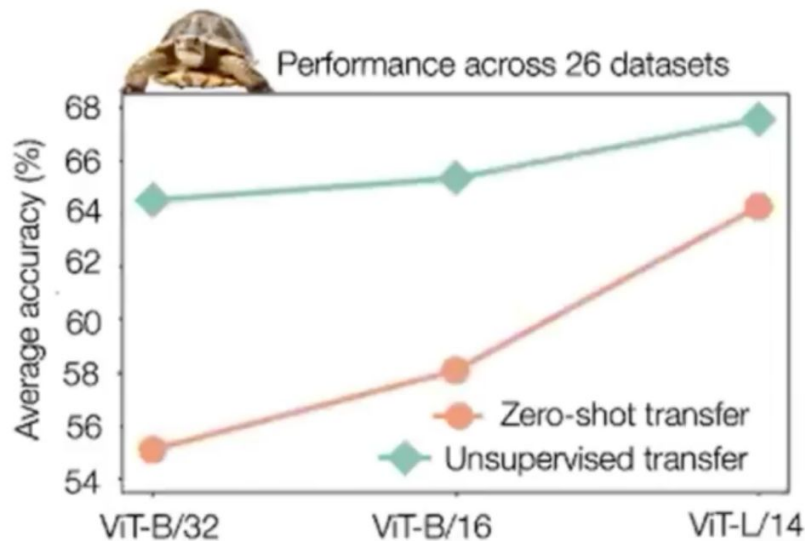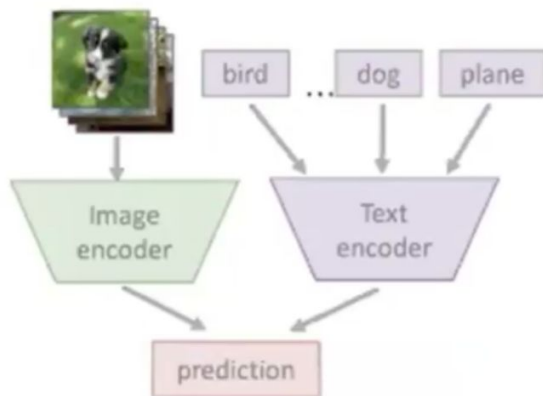Artyom Gadetsky     Yulun Jiang

# Unsupervised Transfer Outperforms Zero-Shot

**SOTA unsupervised performance**

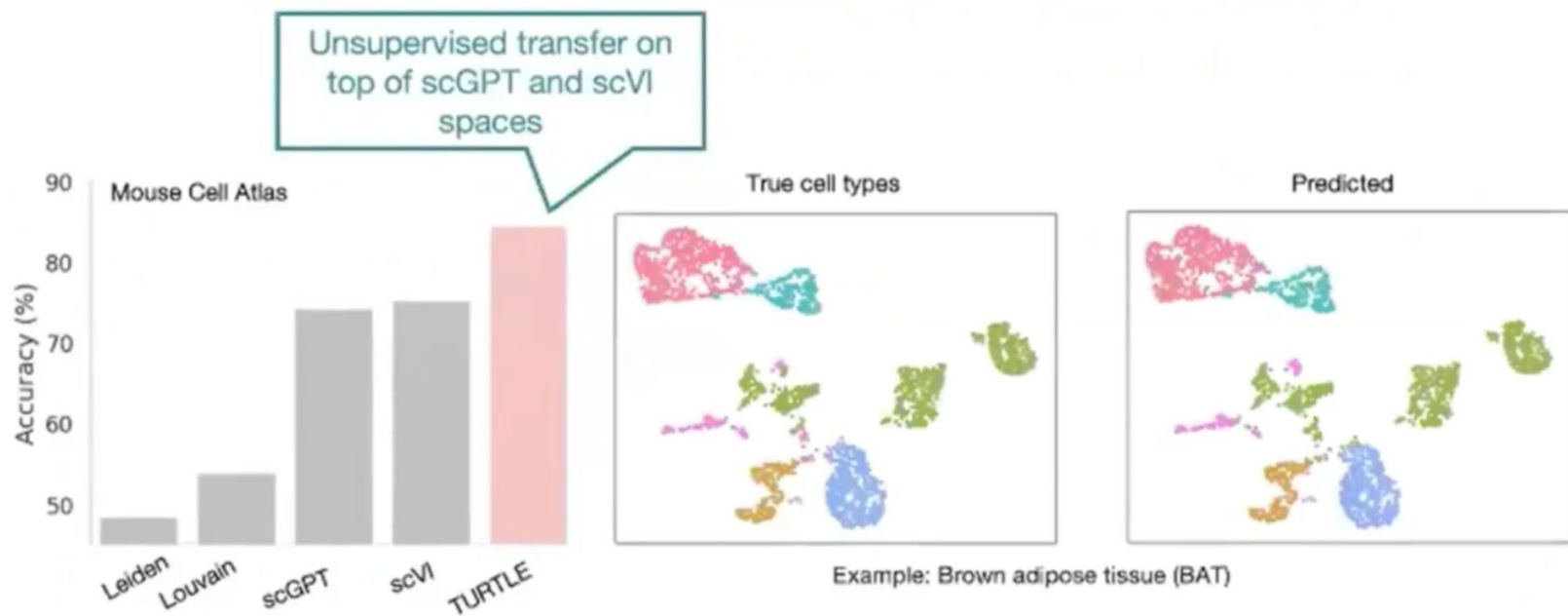- 26 datasets benchmark from CLIP



**OpenAI** Zero-Shot CLIP Model

bird ... dog plane

Image encoder

Text encoder

prediction

Gadetsky*, Jiang*, Brbić. *ICML '24*

Performance across 26 datasets

Zero-shot transfer
Unsupervised transfer

**TURTLE is fully unsupervised!**

# TURTLE's Performance Is Correlated to Linear Probe

On 5 datasets, where linear probe attains near perfect accuracy, TURTLE closely matches it!

TURTLE can infer "optimal" classifier without supervision given high-quality representations

$\rho = 0.87$, $p = 6.3 \times 10^{-9}$

Unsupervised TURTLE Performance (%) — Supervised Linear Probe Performance (%)

STL10, CIFAR10, Flowers102, Food101, EuroSAT, OxfordPets, RESISC45, Caltech101, UCF101, CIFAR100, StanfordCars, MNIST, SUN397, ImageNet, HatefulMemes, DTD, SST2, Birdsnap, PatchCamelyon, GTSRB, KITTI Distance, Kinetics700, FGVCAircraft, FER2013, CLEVRCounts, Country211

Mario Ortai, EPFL

# Application to Single-Cell Data



Unsupervised transfer on top of scGPT and scVI spaces

Mouse Cell Atlas

True cell types

Predicted

Example: Brown adipose tissue (BAT)

# AlphaFold For Cells

Protein structure

Protein sequence

# AlphaFold For Cells



Gene expressions

Genes

Cells

Tissue structure

# LUNA: From Cells to Locations

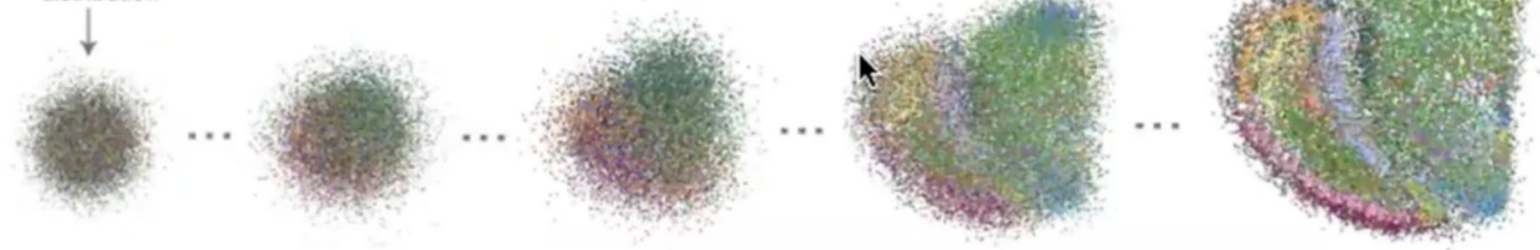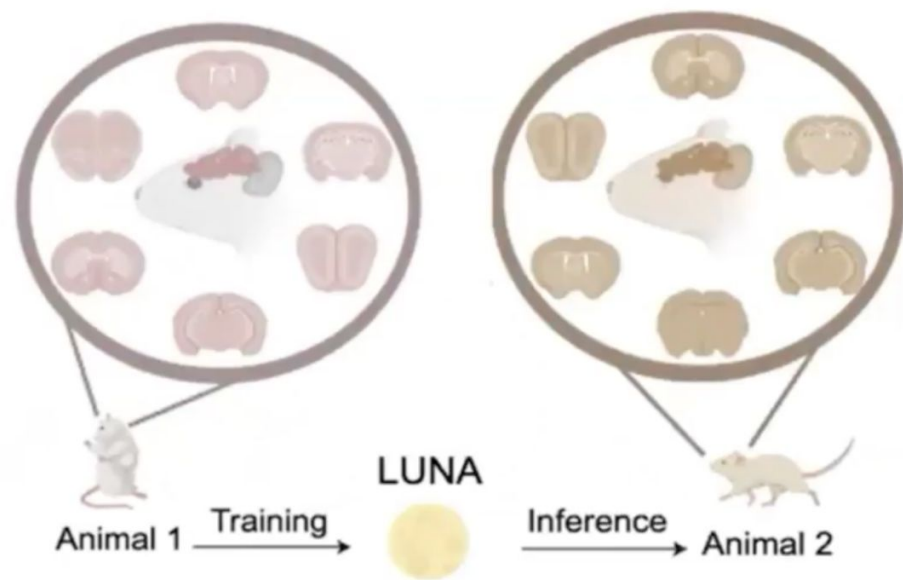A generative model for mapping cells to their locations and generating tissue structures

Yist Yu    Chanakya Ekbote

Sample from standard normal distribution

*Unpublished work*

# Reconstruction of Whole Mouse Brain MERFISH Atlas



LUNA

Animal 1 →Training→ Inference→ Animal 2

*Unpublished work*

Dataset:
- MERFISH Mouse Brain Atlas with over 4 million cells

Training dataset:
- 2.85 million cells across 147 slices from one mouse

Target unlabeled dataset:
- 1.23 million cells across 66 slices from another mouse

33

# Reconstruction of Whole Mouse Brain MERFISH Atlas



Prediction

Groundtruth

*Unpublished work*

# Reconstruction of Whole Mouse Brain MERFISH Atlas

**338 different subclasses!**

Ground truth

LUNA's predictions

# Reconstruction of Whole Mouse Brain MERFISH Atlas

# LUNA: Zero-Shot Setting

Train set

Test set

Animal 1

Animal 2

Cell classes:

Cell classes:

unseen
cell class

# Zero-Shot Generalization to Unseen Cell Types

NP-CT-L6b Glut

Tshz2 gene

Cell class unseen during the model training

Ground truth

LUNA's predictions

*Unpublished work*

39

# De Novo Reconstruction of CNS ScRNA-seq Atlas



Estimated ground truth

LUNA's predictions

216 cell classes

*Unpublished work*

# De Novo Reconstruction of CNS ScRNA-seq Atlas

LPL gene

IQGAP2 gene

MEF2C gene



*Unpublished work*

41

# Acknowledgements



## MLBio lab@EPFL

**PhD students**
**Artyom Gadetsky**
Shuo Wen
**Tingyang Yu**
**Yulun Jiang**
Siba Panigrahi

**Postdocs:**
Ramon Viñas Tórne
Myeongho Jeon

**Research engineers**
Jeremy Goumaz

**Assistant:**
Marie Künzle

Machine Learning

Biomedicine

**Collaborators:**
Jure Leskovec, Stanford
Yanay Rosen, Stanford
Yusuf Roohani, Stanford
Kaidi Cao, Stanford
John Hickey, Duke
Hongjie Li, Baylor College of Medicine
Yuqi Tan, Stanford
Liqun Luo, Stanford
Michael Snyder, Stanford
Garry Nolan, Stanford
Pascal Frossard, EPFL
Chanakya Ekbote, MIT

**Swiss National Science Foundation**

NIH ImmGen

**Filipi Silva,** *Indiana University*

AI / neural networks are complex

**data** → **AI** → **output**

doc2vec, transformers, node2vec, etc
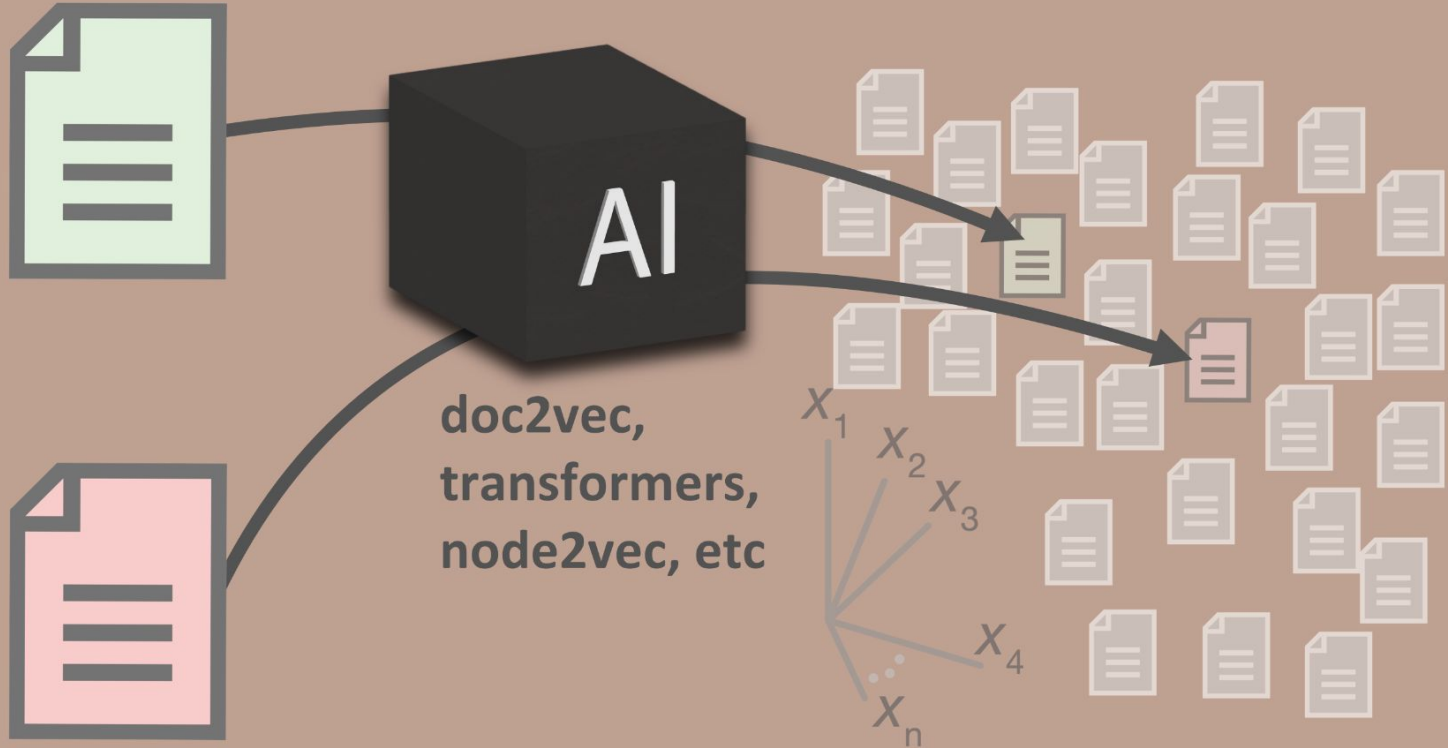
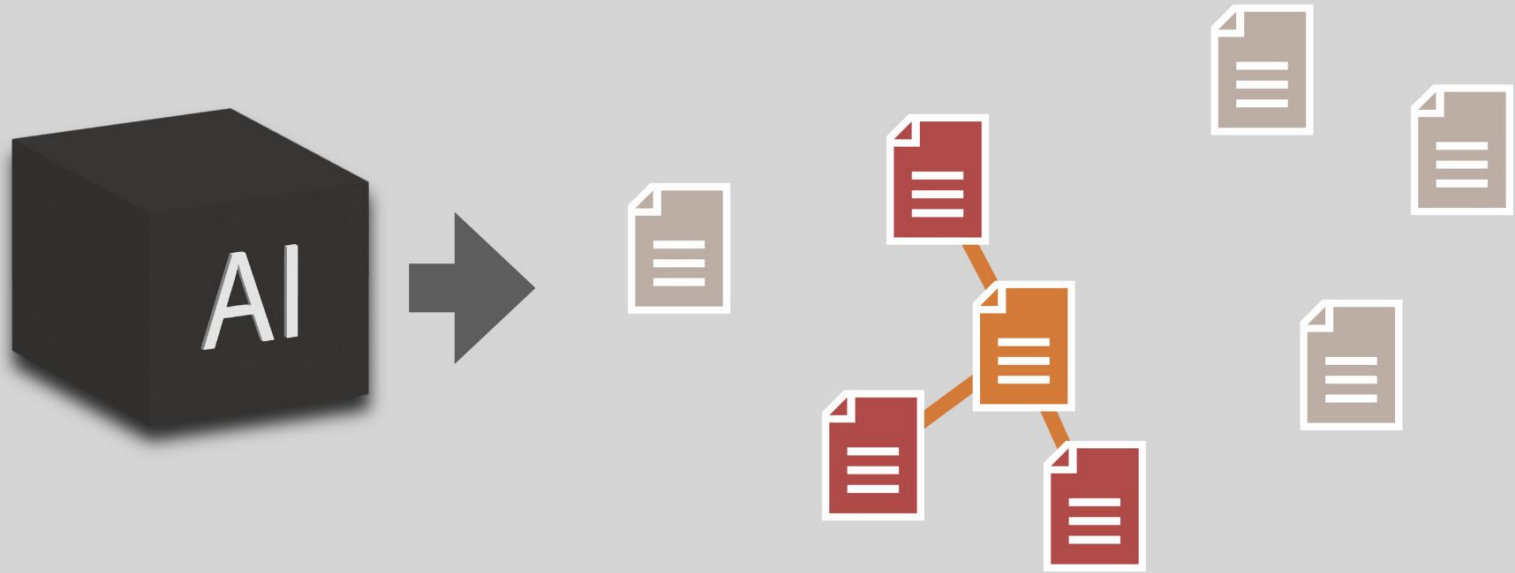$X_1$ $X_2$ $X_3$ $X_4$ $X_n$
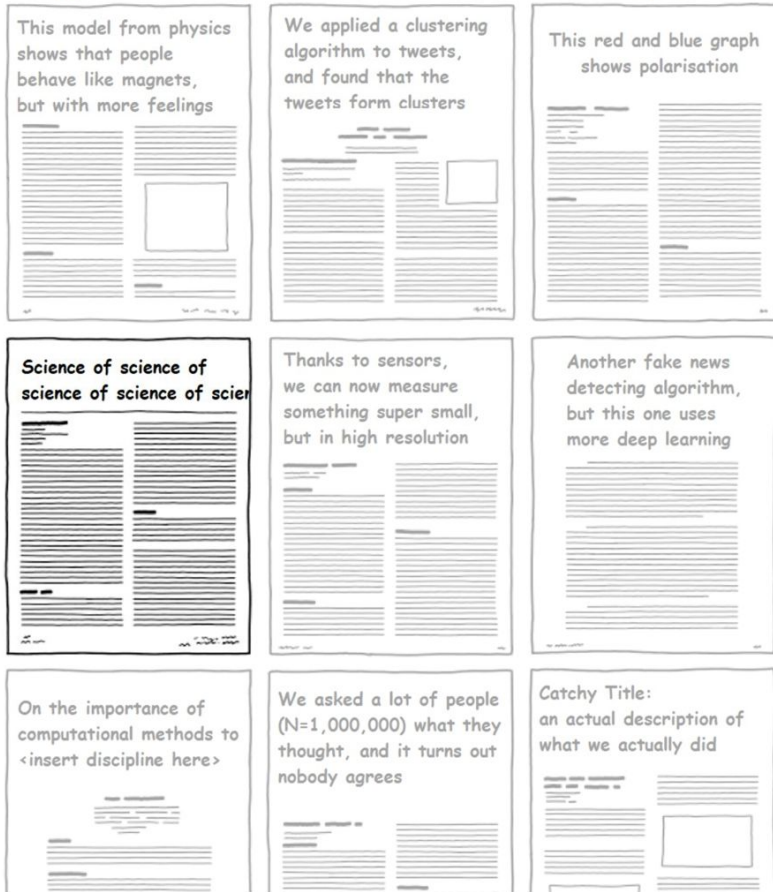
**Embeddings**

**Recommendation systems**

**Databases/Search engine**

**Retrieval-Augmented Generation (RAG)**

**Anomaly Detection**

...

Data is also complex

## Types of Computational Social Science papers

This model from physics shows that people behave like magnets, but with more feelings

We applied a clustering algorithm to tweets, and found that the tweets form clusters

This red and blue graph shows polarisation

Science of science of science of science of scier

Thanks to sensors, we can now measure something super small, but in high resolution

Another fake news detecting algorithm, but this one uses more deep learning

On the importance of computational methods to <insert discipline here>

We asked a lot of people (N=1,000,000) what they thought, and it turns out nobody agrees

Catchy Title: an actual description of what we actually did

by Chico Camargo (Twitter)
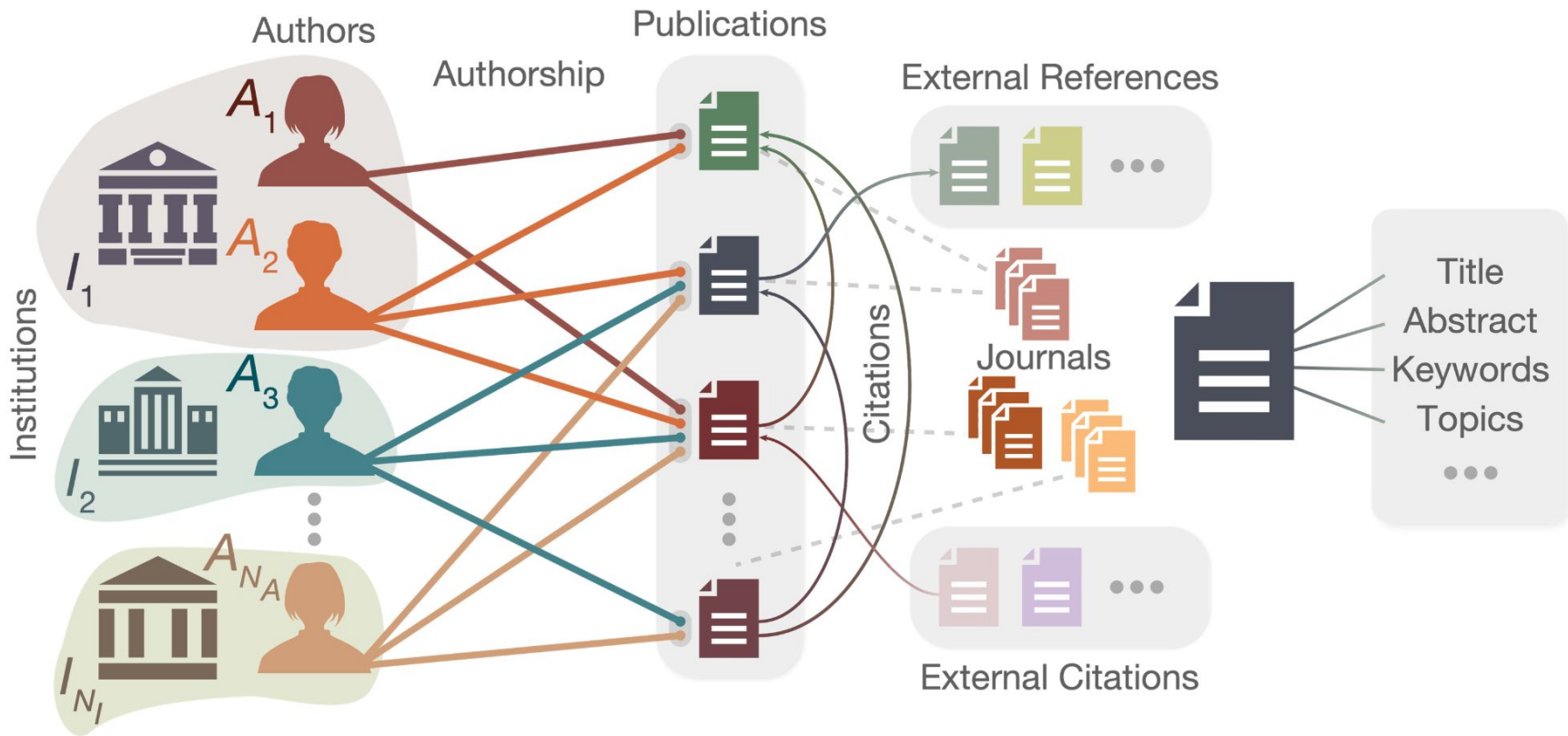https://twitter.com/evoluchico/status/1388137531552718860

# Science of Science

- How science is evolving?

- How researcher teams are formed?

- Is science becoming more interdisciplinary?

- Can we predict success in science?

- How to properly evaluate researchers? journals? papers?

- Can tools/approaches accelerate the scientific development?

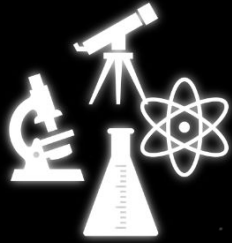- Can we predict the benefits of implementing a policy?
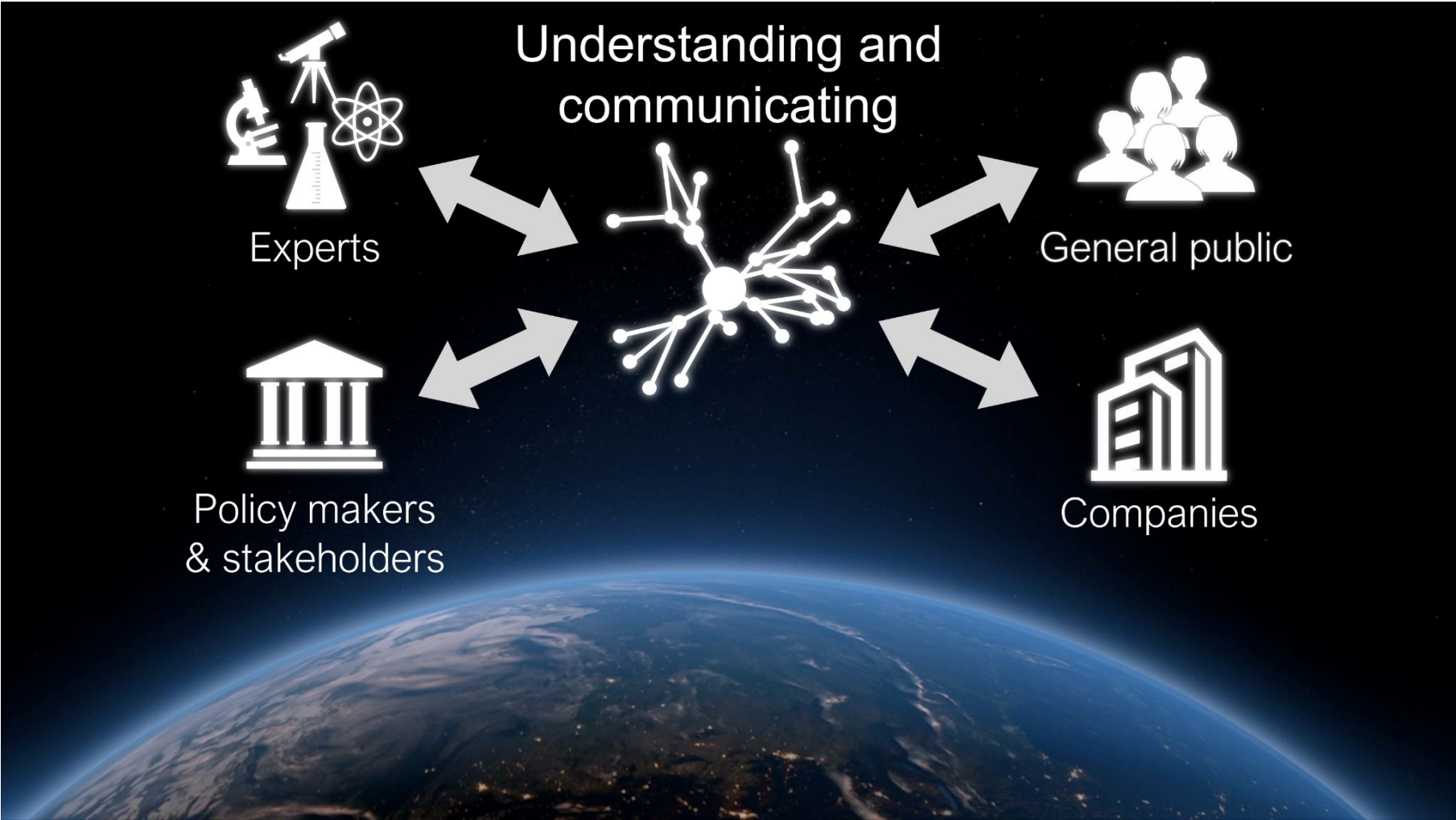
- ...

# Publications

Institutions

Authors

$A_1$

$A_2$

$A_3$

$A_{N_A}$

$I_1$

$I_2$

$I_{N_I}$

Authorship

Publications

Citations

External References

Journals

External Citations

Title
Abstract
Keywords
Topics
· · ·

# Maps

G

Oregon U.
LIGO Hanford Observ.
Birmingham U.
Potsdam, Max Planck Inst.
Syracuse U.
Cardiff U.
Michigan U.
Wisconsin U., Milwaukee
Columbia U.
Caltech
Northwestern U.
NASA, Goddard
LIGO Livingston Obs.
Florida U.

Australian Natl. U., Canberra

● Research Center
● University
● College
● Other

Participating countries over time

| | Co | R | U | R |

80
60
40
20
0

H

2000  2005  2010  2015  2020
Year of first collaboration

LIGO Hanford Observ.
Oregon U.
Wisconsin U., Milwaukee
Syracuse U.
Northwestern U.
Michigan U.
Columbia U.
Caltech
NASA, Goddard
LIGO Livingston Obs.
Florida U.

I

Börner, K., Silva, F. N., & Milojević, S. (2021).
Visualizing big science projects. Nature Reviews Physics, 3(11), 753-761.

21

Leydesdorff, Loet, and Ismael Rafols. "A global map of science based on the ISI subject categories." Journal of the American Society for Information Science and Technology 60.2 (2009): 348-362.

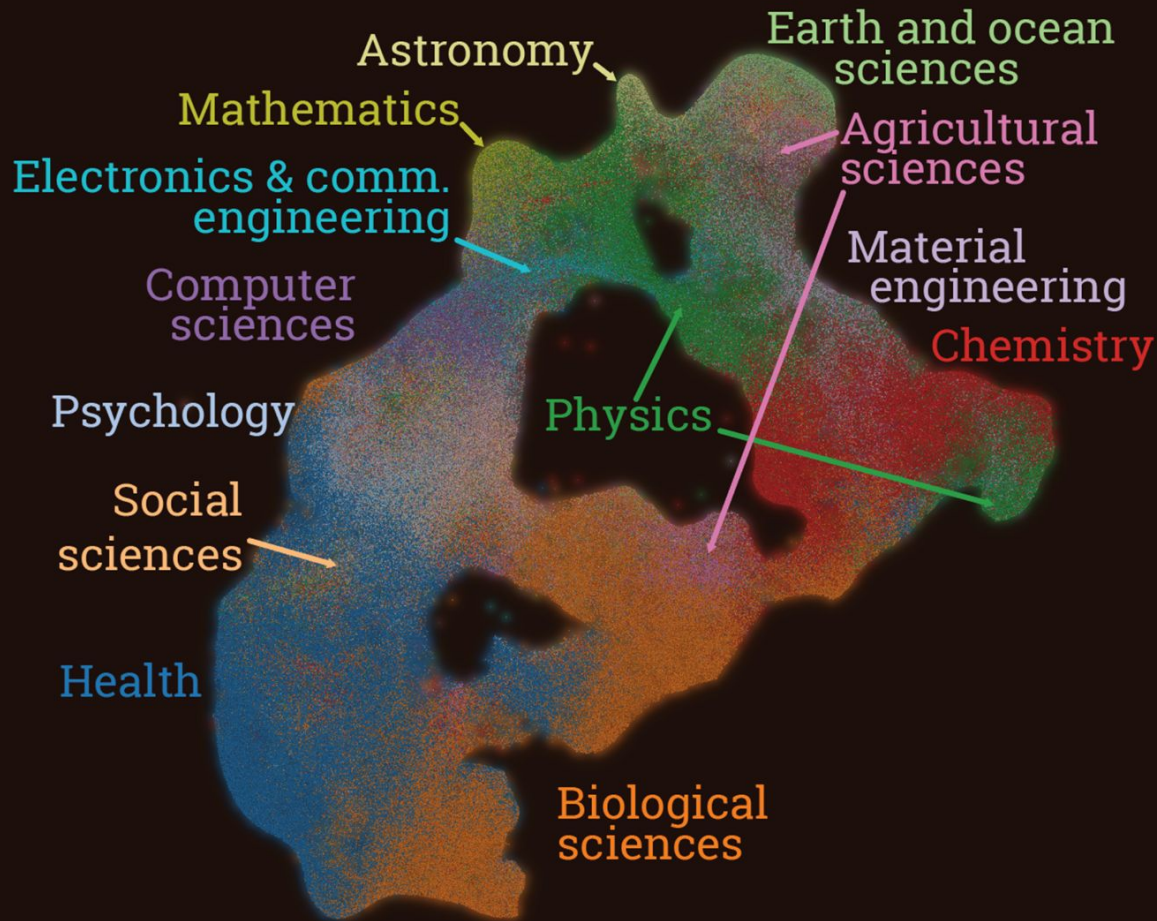ETO Map of Science
sciencemap.eto.tech

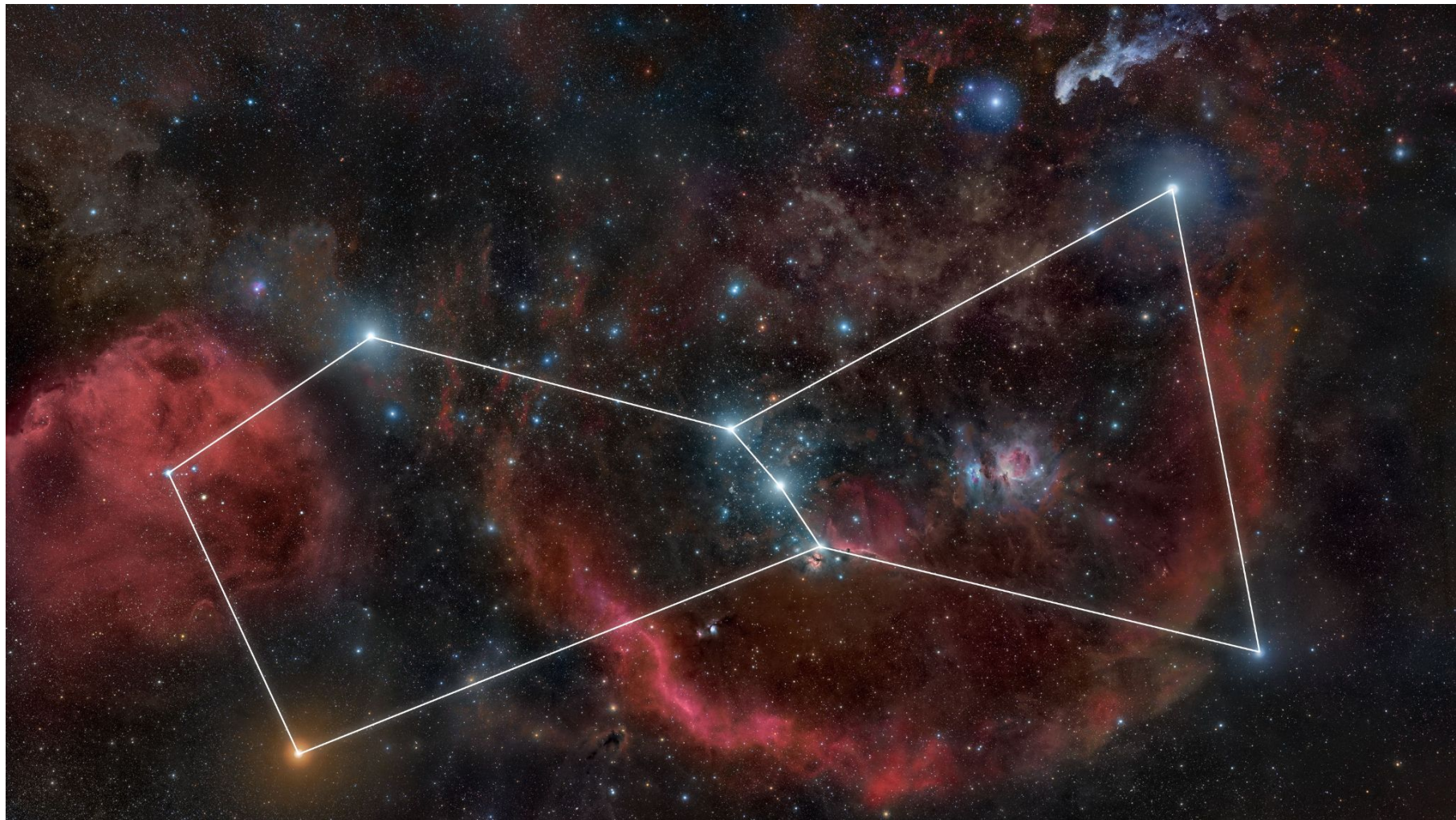Bollen, Johan, et al. "Clickstream data yields high-resolution maps of science." PloS one 4.3 (2009): e4803.

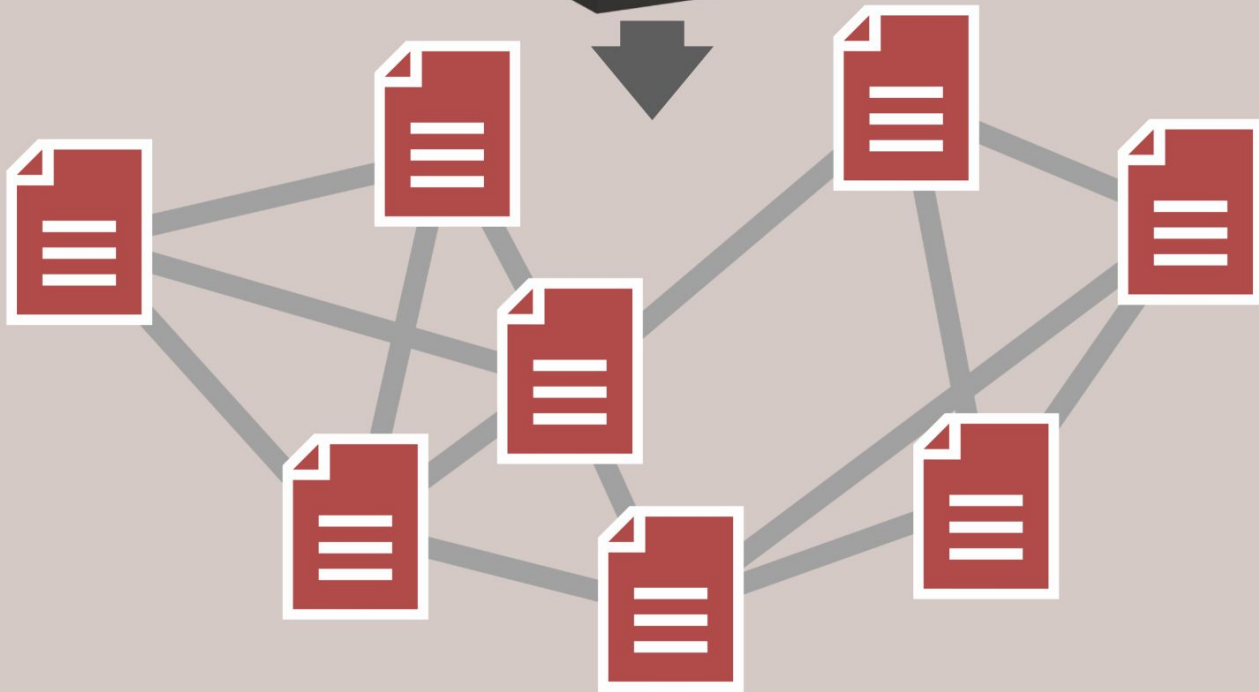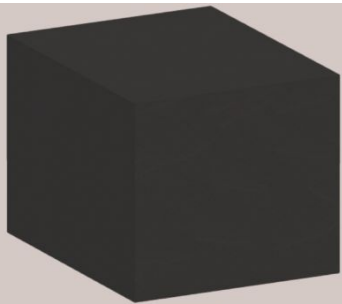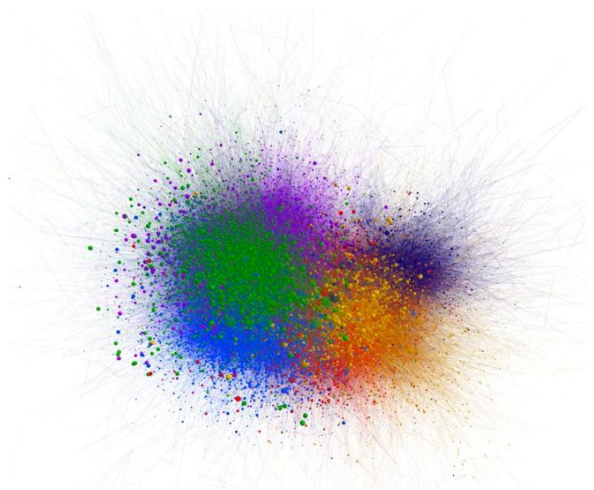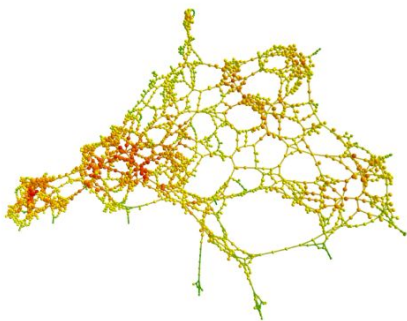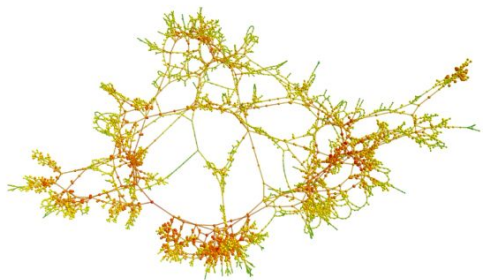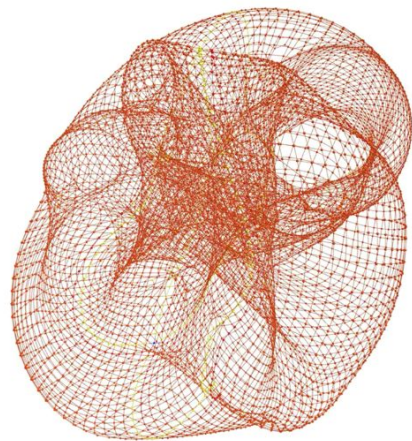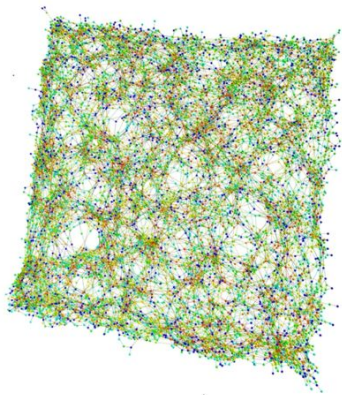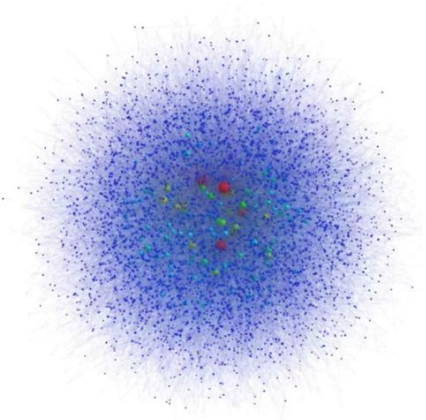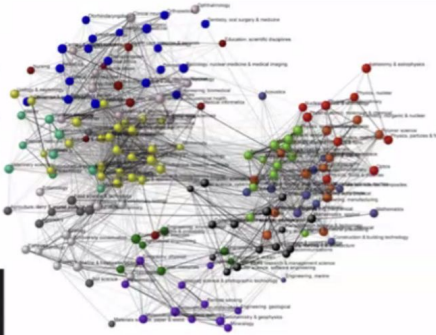Börner, Katy, et al. "Design and update of a classification system: The UCSD map of science." *PloS one* 7.7 (2012): e39464.
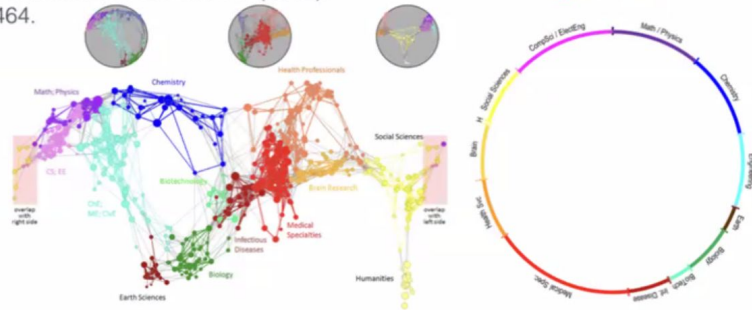
**data**

Leydesdorff, Loet, and Ismael Rafols. "A global map of science based on the ISI subject categories." Journal of the American Society for Information Science and Technology 60.2 (2009): 348-362.

ETO Map of Science
sciencemap.eto.tech

Börner, Katy, et al. "Design and update of a classification system: The UCSD map of science." *PloS one* 7.7 (2012): e39464.
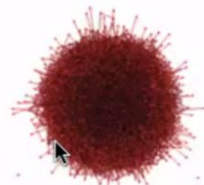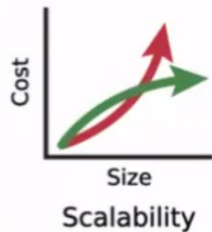
Bollen, Johan, et al. "Clickstream data yields high-resolution maps of science." PloS one 4.3 (2009): e4803.

**Network characteristics**

Large-scale

Dynamic

Complex

Scalability

Overcome "Hairballs"

Capture the multiscale structure

Real-time dynamic exploration

**M1 - Layout algorithms**

Continuous layout via sparse-matrix operations

Multi-scale trajectories representation

**M2 - Rendering pipeline**

SDF, billboards GPU-based

Edge density via adv. blending

**M3 - Interactivity**

Interactive network transformations

Multiple representations

# web Helios

**Open-source web framework**
can be integrated in websites, portals, dashboards ...

**Optimized rendering and layouts**
can visualize large networks, high-quality rendering ...

**Interactivity***
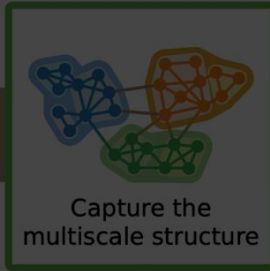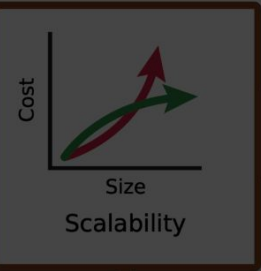allows picking, filtering, navigation, multi-representations ...

*in development

$F_r$

$F_a + F_r$

Repulsive Force

Attractive Force

# Layout optimizations

- Molecular dynamics simulation is $O(N^2)$ .

- We can use multipole expansion (FM3):

- Segment the space

- Real-time continuous layout



* S. Hachul and M. Jünger. Drawing large graphs with a potential-field-based multilevel algorithm. In International Symposium on Graph Drawing, pp. 285–295. Springer, 2004. doi: 10.1007/978-3-540-31843-9_29

# Wiki Medicine and Mathematics

# Rendering in the GPU

**a** Billboards for nodes and edges

Edge

Node

**b** Rendering shapes

Circle

2D SDF

3D SDF

Texture

render

render encoding indices

Camera

**c** View and picking framebuffers

view buffer

2
4
3
1

picking buffer

decode pixel

pick event node/edge

fury.gl

Demo made by Javier Guaje  Serge Koudoro and Eleftherios Garyfallidis

Rendering

**Edge density**

**Ambience Occlusion**

mass, neutrino, standard model, decay, boson, higgs

atom, trap, bose-einstein condensate, gas, optical, interaction

electron, ion, calculation, cross section, ionization, energy, collision'

nucleus, reaction, mev, nuclear, energy, neutron, calculation'

quantum, state, entanglement system, show, photon

spin, temperature, superconducting, electron, magnetic field

field, cosmological, theory, scalar, universe, gravity

decay, meson, quark, qcd, pi

liquid, simulation, dynamic, fluid, surface

system, model, dynamic, network, time

A

A - artificial neural network, simulation, system, pape...
B - convolutional neural network, propose, learning, t...
C - kalman filter, estimation, estimate, state, measur...
D - support vector machine, classification, classifier, ...
E - signal, brain, independent component analysis, e...
F - image, patient, classification, diagnosis, %, disea...
G - hidden markov model, feature, paper, speech rec...
H - linear regression, study, conclusion, associate, a...
I - kriging, spatial, soil, study, sample, area, concentr...
J - principal component analysis, sample, componen...
K - principal component analysis, image, feature, faci...
L - image, classification, area, spatial, spectral, data, ...
M - neural network, invention disclose, accord, meth...
N - neural network, neuron, spike, synaptic, brain, ac...
O - probability, probabilistic, inference, model, struct...
P - fault diagnosis, base, monitoring, propose, signal, ...
Q - molecular, compound, descriptor, regression, rel...
R - linear regression model, estimator, estimate, nonl...
Other

Export   Size ▢▢   Color  cluster name (level1   Category18     Edges ▮▢

**Ok, but what about embeddings?**

UMAP uses SGD but could use other FD

$F_r$
$F_a + F_r$

Repulsive forces
Attractive forces

For large embeddings: >10M points

**Negative samples rate :** Number of repulsive interactions to update for each positive.

increased to 10 (default = 5)

**Epochs:** Number of iterations

increased to 200000 (default = 200) !!!

**Number of neighbors:** Number of neighbors in the NN graph

increased to 30 (default = 15) !!!

umap-learn.readthedocs.io

attractive        repulsive

$$\frac{\partial \mathcal{L}_{\text{UMAP}}(\gamma)}{\partial \mathbf{y}_i} \sim \sum_j v_{ij} w_{ij}(\mathbf{y}_i - \mathbf{y}_j) - \gamma \sum_j \frac{1}{d_{ij}^2 + \epsilon} w_{ij}(\mathbf{y}_i - \mathbf{y}_j).$$

**GPU version (much faster!)**

beware: it has bugs that lead to bad projections
https://docs.rapids.ai/api/cuml/stable/api/#umap

Böhm, J. N., Berens, P., & Kobak, D. (2022). Attraction-repulsion spectrum in neighbor embeddings. The Journal of Machine Learning Research, 23(1), 4118-4149.

Astronomy

Mathematics

Earth and ocean sciences

Electronics & comm. engineering

Agricultural sciences

Computer sciences

Material engineering

Chemistry

Psychology

Physics

Social sciences

Health

Biological sciences

Interactive version

Computer sciences

Psychology+brain interdisciplinary

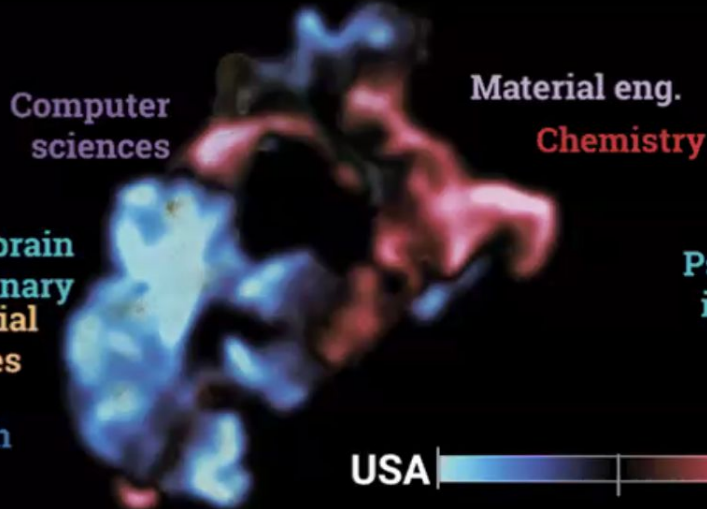General & internal medicine

Health

Cancer

RNA extraction

Genomics

Citation Count

Where there is disruption?

Computer sciences

Phys., ag., eng., mat. sciences interface

Psychology+brain interdisciplinary

Agricultural sciences

Top 5% disruptive Number of papers

Who produces what?

Computer sciences

Material eng. Chemistry

Psychology+brain interdisciplinary Social Sciences

Health

USA ⟷ China

Who funds disruptive science?

Computer sciences

Phys., ag., eng., mat. sciences interface

Psychology+brain interdisciplinary

Agricultural sciences

A - Recent publications (2010 onwards) by authors who coauthored COVID-related papers

Materials sci. eng.

Chemistry

Psychology    Health

Physics

Chemical eng.

Astronomy

Social sci.

Electrical eng.    Agricultural

Earth and ocean

Mathematics

Biological sci.

Computer sci.    Atmospheric

Mechanical eng.

Other eng.

B - COVID-related papers

Materials sci. eng.

Physics    Chemistry    Psychology    Health

Astronomy    Chemical eng.

Social sci.

Electrical eng.    Agricultural

Mathematics    Earth and ocean

Computer sci.    Biological sci.

Atmospheric

Mechanical eng.
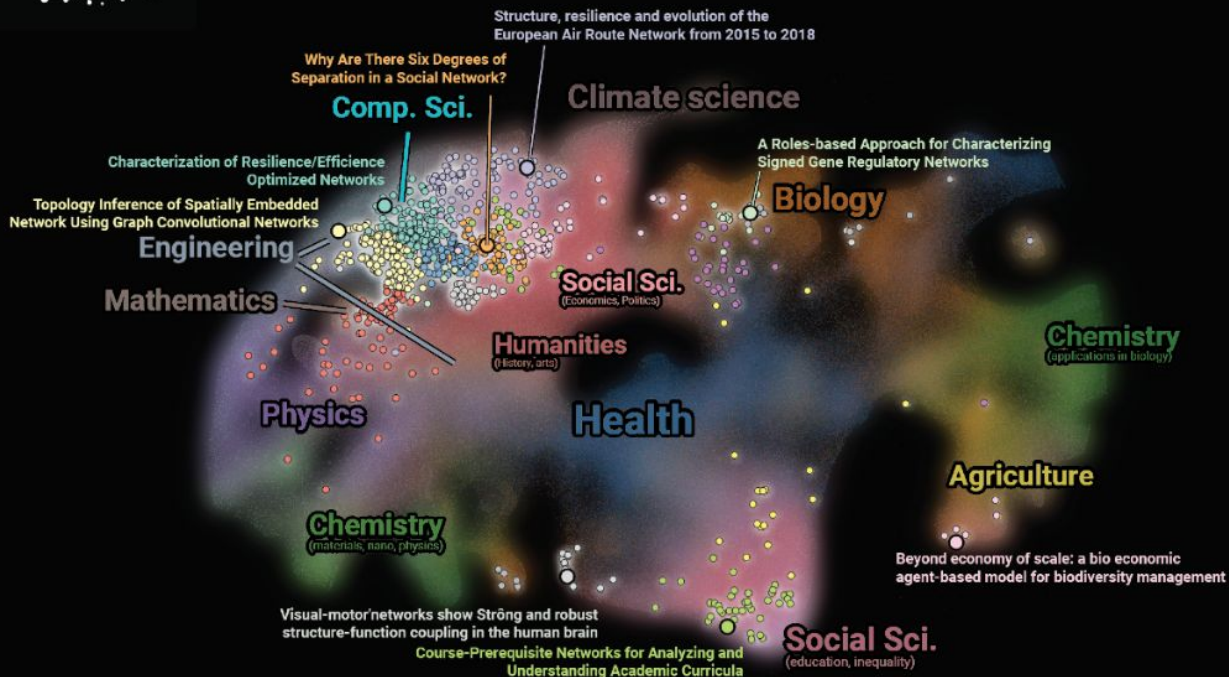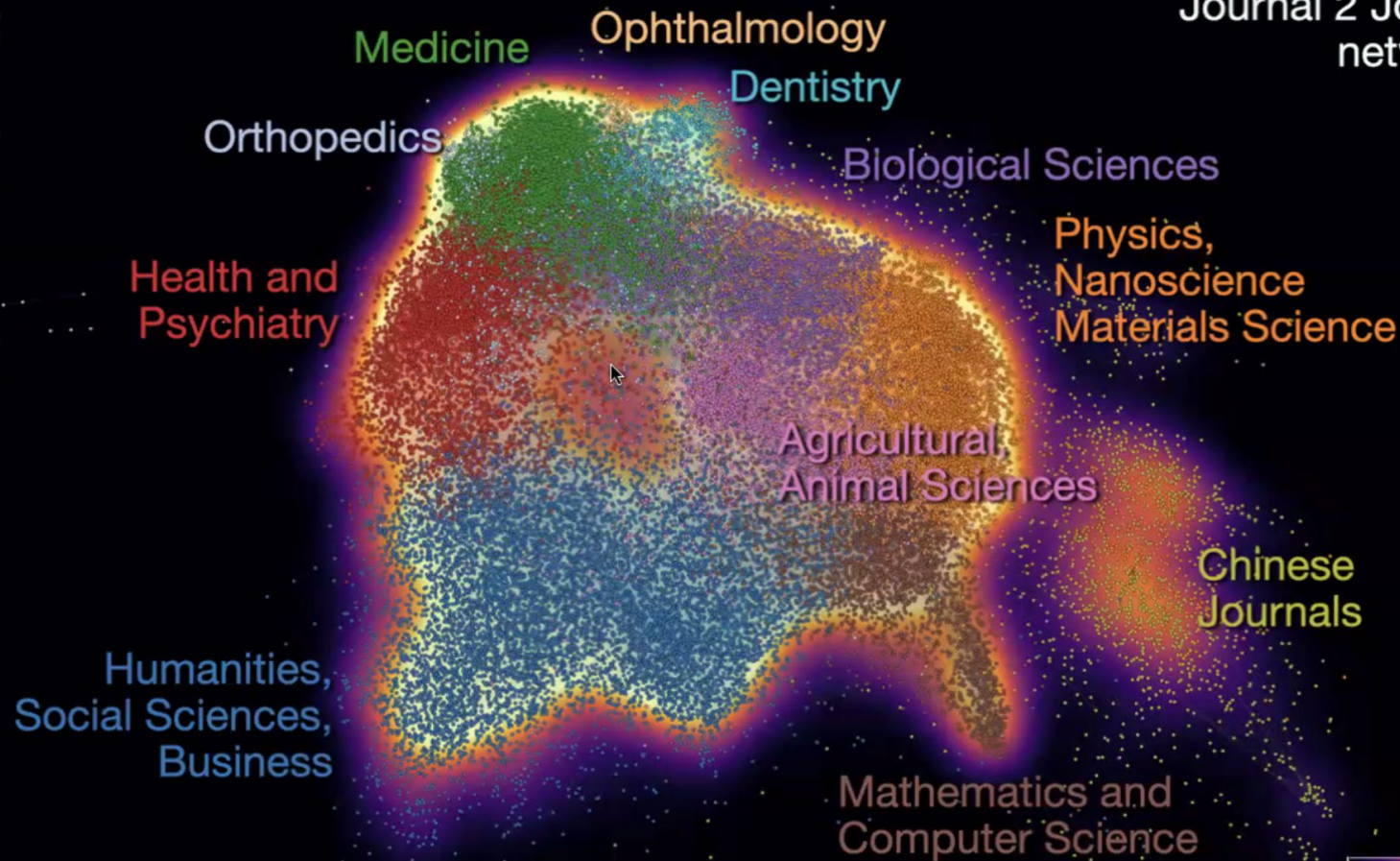
Other eng.

Where in the world of knowledge is NetSci?

A science map constructed from the titles and abstracts of publications in the Web of Science. On top of that, we project NetSci contributions from 2023 and 2024.

Atmospheric sciences
Earth and ocean sciences
Chemistry
Materials science engineering
Physics
Astronomy
Chemical engineering
Agricultural sciences
Other engineering
Mechanical engineering
Civil engineering
Electrical / electronics and communications engineering
Computer sciences
Biological sciences
Social sciences
Psychology
Health

Psychology
Health
Biological sciences
Social sciences
Agricultural sciences
Computer sciences
Atmospheric sciences
Earth and ocean sciences
trical, electronics and communications engineering
Other engineering
Chemical engineering
Mathematics
MechCivil engineering.g
Chemistry
Materials science engineering
Physics
Astronomy

Node2vec version

SPECTER (titles + fine tuning via citation network)
embedding of the whole science
Microsoft Academic Graph (more than 200M works)

Ongoing project with YY Ahn, Sadamori Kojaku and others

$$\text{sim}(a, b) \propto \log\left[\frac{P_{a \to b}}{P_a P_b}\right]$$

**Map for Physics**

Particles & Fields

Astronomy & Astrophysics

Nuclear Physics

**1**

Mechanics

Interdisciplinary Physics

Fluid & Plasmas

**2**

Optics

Atomic molecular & Chemical

Applied Physics

Condensed Matter

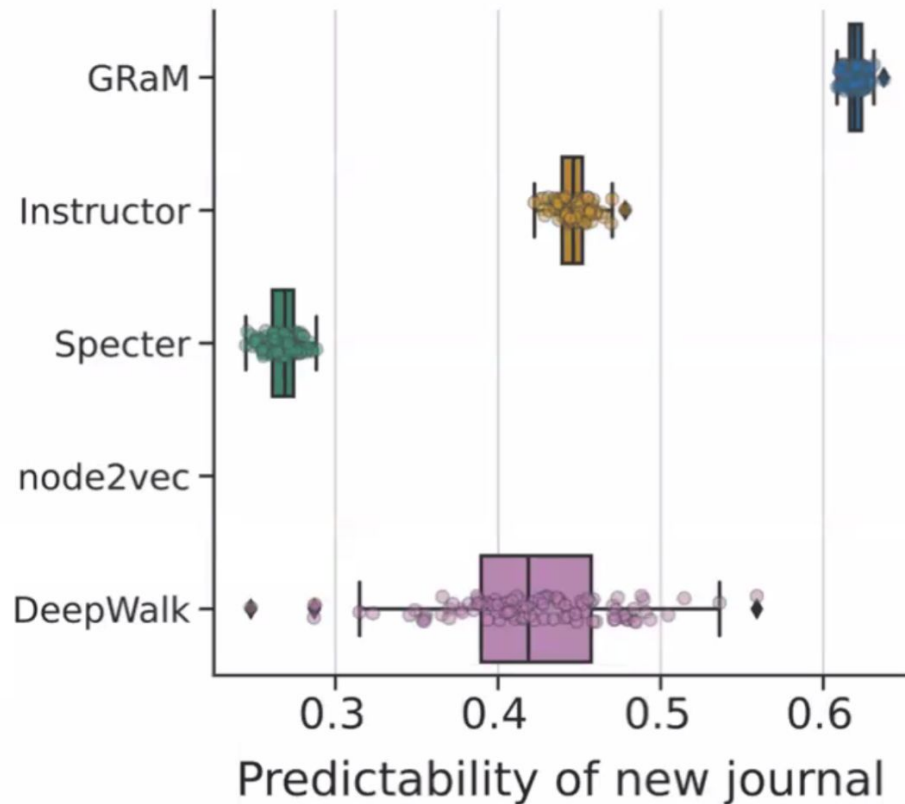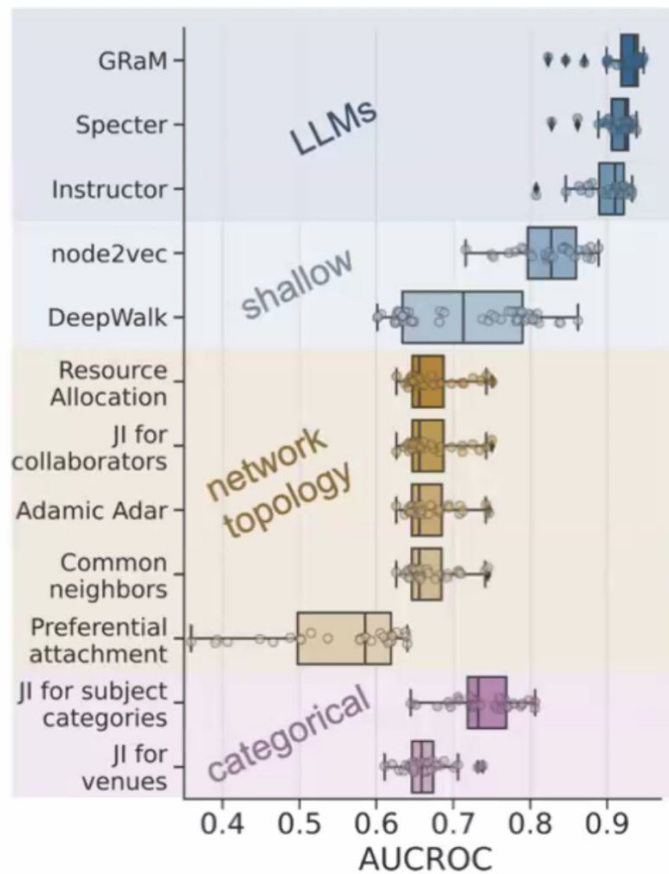**Trajectories of Nobel laureates**

Giorgio Parisi

Donna Strickland

**How well it represents topics?**

Do neighboring papers share the same subject categories?

Instructor

SPECTER

GRaM

node2vec

DeepWalk

In/Out Citations

Precision

Number of neighbors, $k$

Predictive power of the GRaM

http://osome.iu.edu/tools/networks

https://osome.iu.edu/tools/coordiscope

# Adapting these models and tools for biomedical research



- Single cell RNA-seq data

- Microbiome associations with tissue gene expression

- Brain networks across age

- Metabolomics

# Cell Type Differentiation Using Community Detection



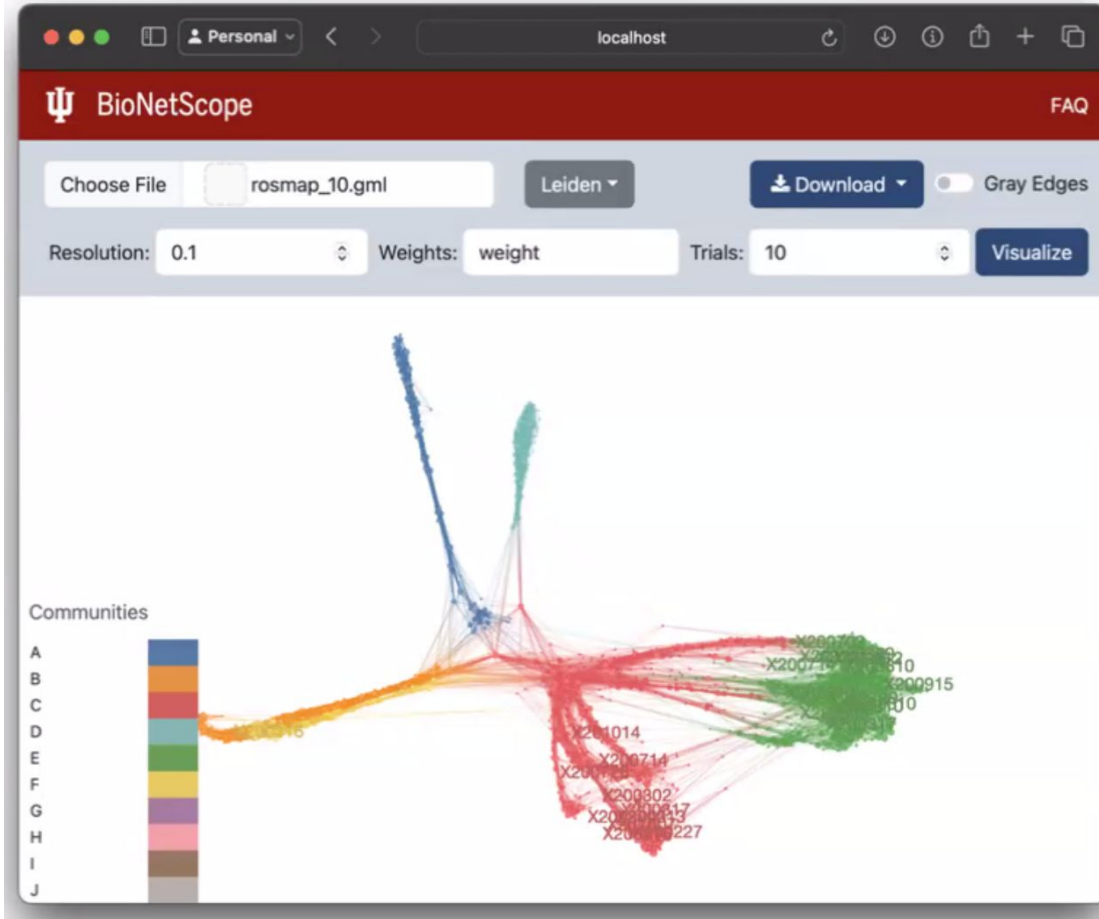68k human Peripheral Blood Mononuclear Cells (PBMCs) scRNA-seq dataset

Fatemi Nasrollahi, F. S., Silva, F. N., Liu, S., Chaudhuri, S., Yu, M., Wang, J., ... & Fortunato, S. (2024). Cell Type Differentiation Using Network Clustering Algorithms. bioRxiv, 2024-12.

Legend:
- Myocyte (sk. muscle)
- Fibroblast I
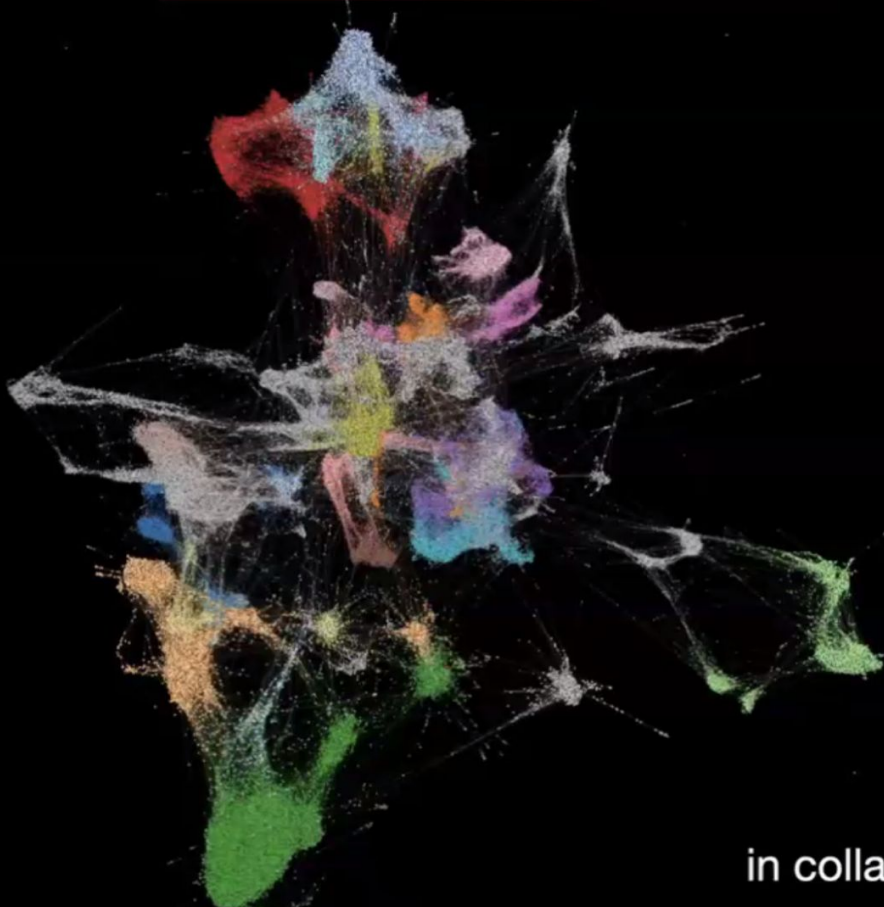- Epithelial cell (alveolar type II)
- Epithelial cell (luminal)
- Endothelial cell (vascular) I
- Myocyte (smooth muscle TAGLN lo)
- Fibroblast
- Epithelial cell (basal I)
- Endothelial cell (vascular) II
- Epithelial cell (basal II)
- Epithelial cell (alveolar type I)
- Myocyte (cardiac)
- Myocyte (smooth muscle)
- Fibroblast II
- Immune (macrophage I)
- Endothelial cell (lymphatic)
- Immune (alveolar macrophage)
- Epithelial cell (club)
- Other

Search

?

press space to start the layout

SVG   Size ▢▢   Color Granularcelltype   Category18   Edges ▢▢

in collaboration with Katy and others

# Thanks

## Indiana University U.S.

Staša Milojević
Yong-Yeol "YY" Ahn
Santo Fortunato
Filippo Menczer
Alessandro Flammini
Attila Varga
Lili Miao
Katy Börner
Andy Saykin
Thomas M O'Connell
Vijay R. Ramakrishnan
Filippo Radicchi

## George Mason US

Henrique F. de Arruda
Sandro Reia

## Univ. São Paulo Brazil

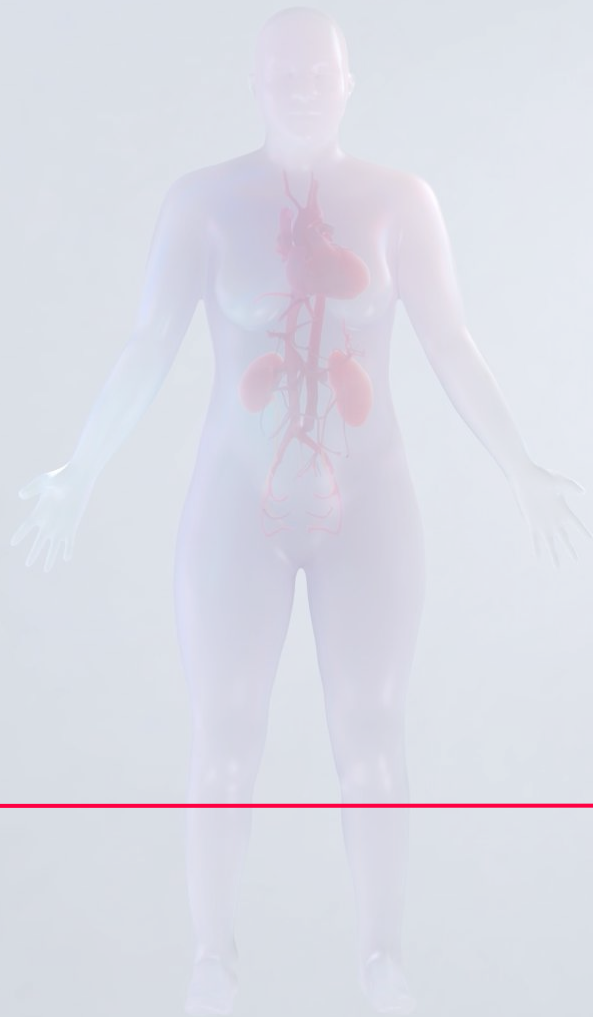Diego R. Amancio
Osvaldo N. de Oliveira Jr.
Ana C. Medeiros

## Binghamton U. U.S.

Sadamori Kojaku

DARPA

NSF

AIR FORCE OFFICE OF SCIENTIFIC RESEARCH — UNITED STATES AIR FORCE

# Q&A

# Questions

How do we define a Multiscale Human?

How do we map a Multiscale Human?

How do we model a Multiscale Human?

How can LLMs or RAGs be used to advance science and clinical practice?

# Thank you