

Data integration and multi-omics

Thomas Stoeger
thomas.stoeger@northwestern.edu

November 24th, 2020

How this lecture will work:

How this lecture will work:

Q/A on handout

How this lecture will work:

Q/A on handout

mini-presentation

How this lecture will work:

Q/A on handout

mini-presentation

group coding exercise

How this lecture will work:

Q/A on handout

**get overview
know where to find details**

mini-presentation

group coding exercise

How this lecture will work:

Q/A on handout

**get overview
know where to find details**

mini-presentation

solidify main points

group coding exercise

How this lecture will work:

Q/A on handout

**get overview
know where to find details**

mini-presentation

solidify main points

group coding exercise

experience

How this lecture will work:

Q/A on handout

**get overview
know where to find details**

mini-presentation

solidify main points

group coding exercise

experience

Q/A on handout

**get overview
know where to find details**

Q/A on handout

get overview
know where to find details

Handout (7min to read)

https://github.com/tstoeger/course_multi_omics/blob/main/multi-omics.pdf

Q/A on handout

get overview
know where to find details

Handout (7min to read)

https://github.com/tstoeger/course_multi_omics/blob/main/multi-omics.pdf

Questions (~5min)

Ask in person.

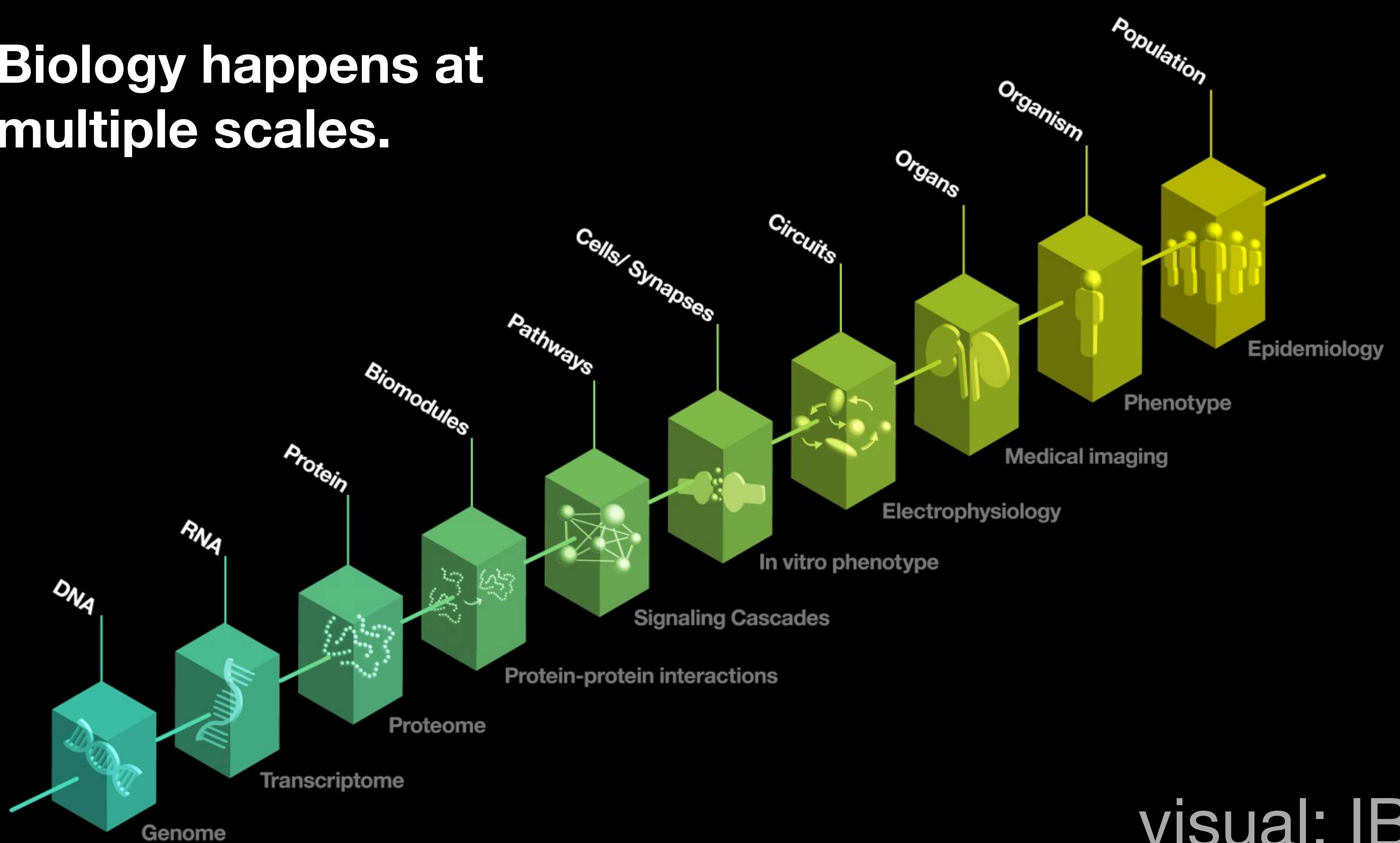
Ask anonymously on:

<https://padlet.com/thomasstoeger/a4mtvpym671dwgnl>

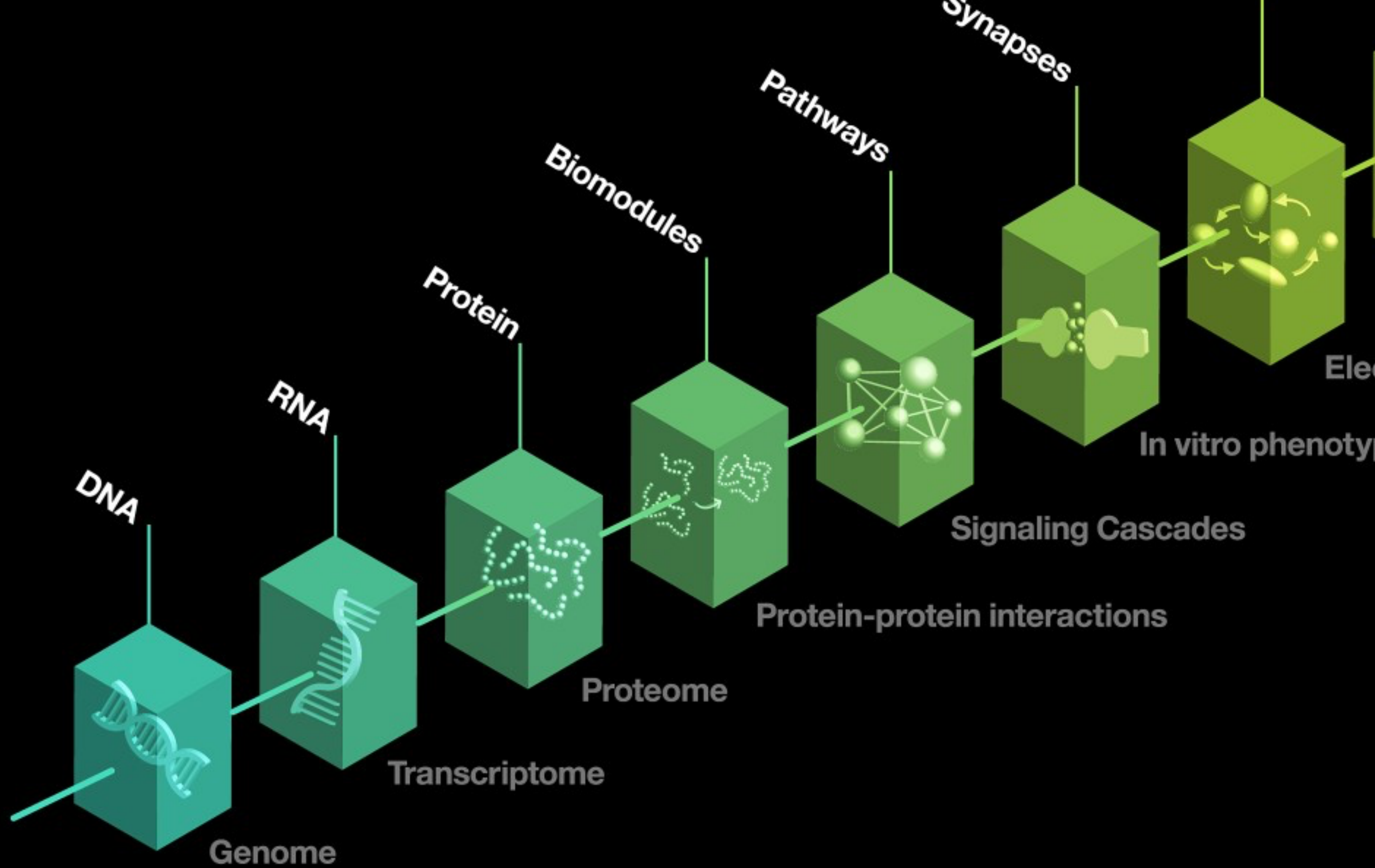
mini-presentation

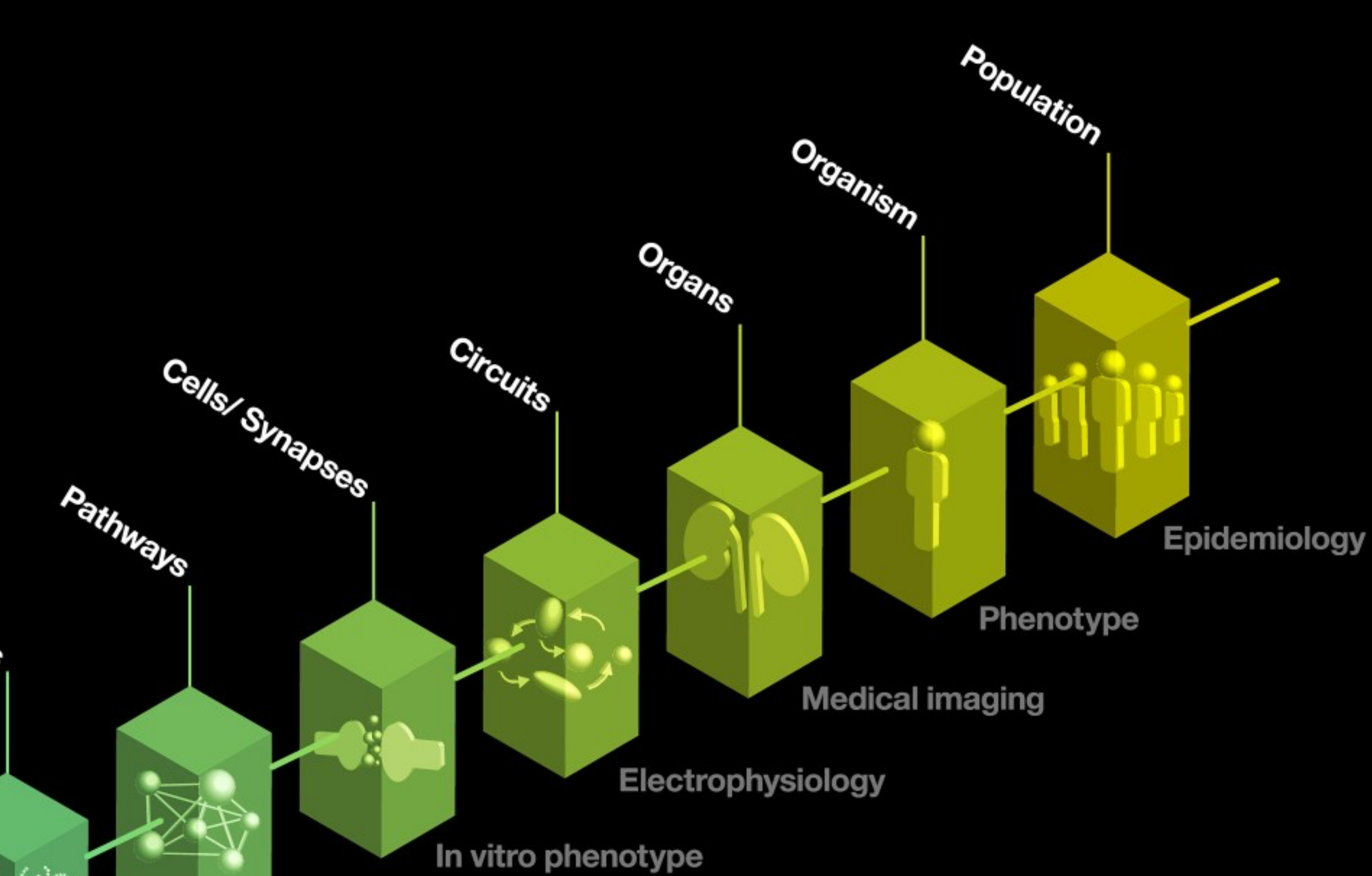
solidify main points

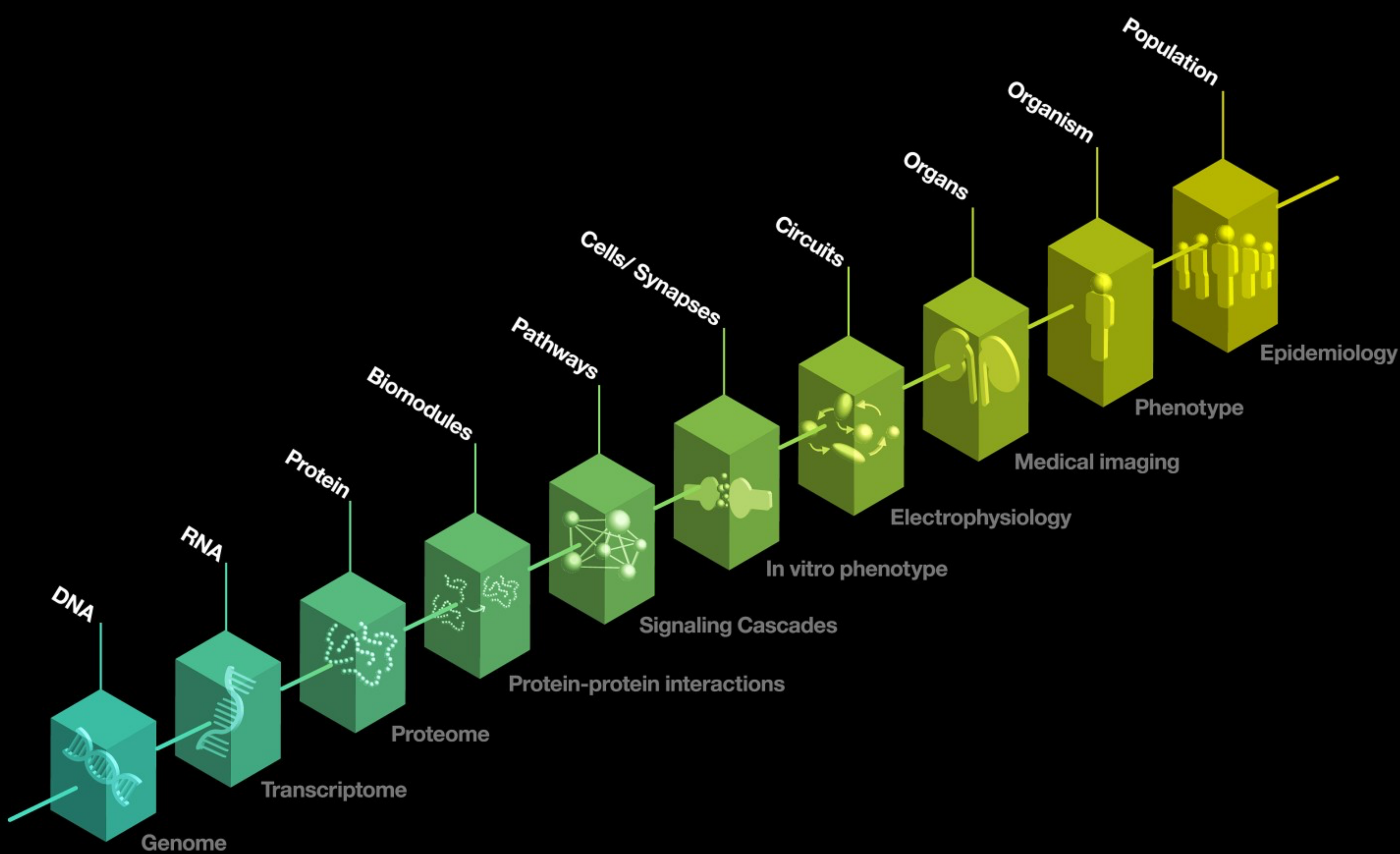
Biology happens at multiple scales.



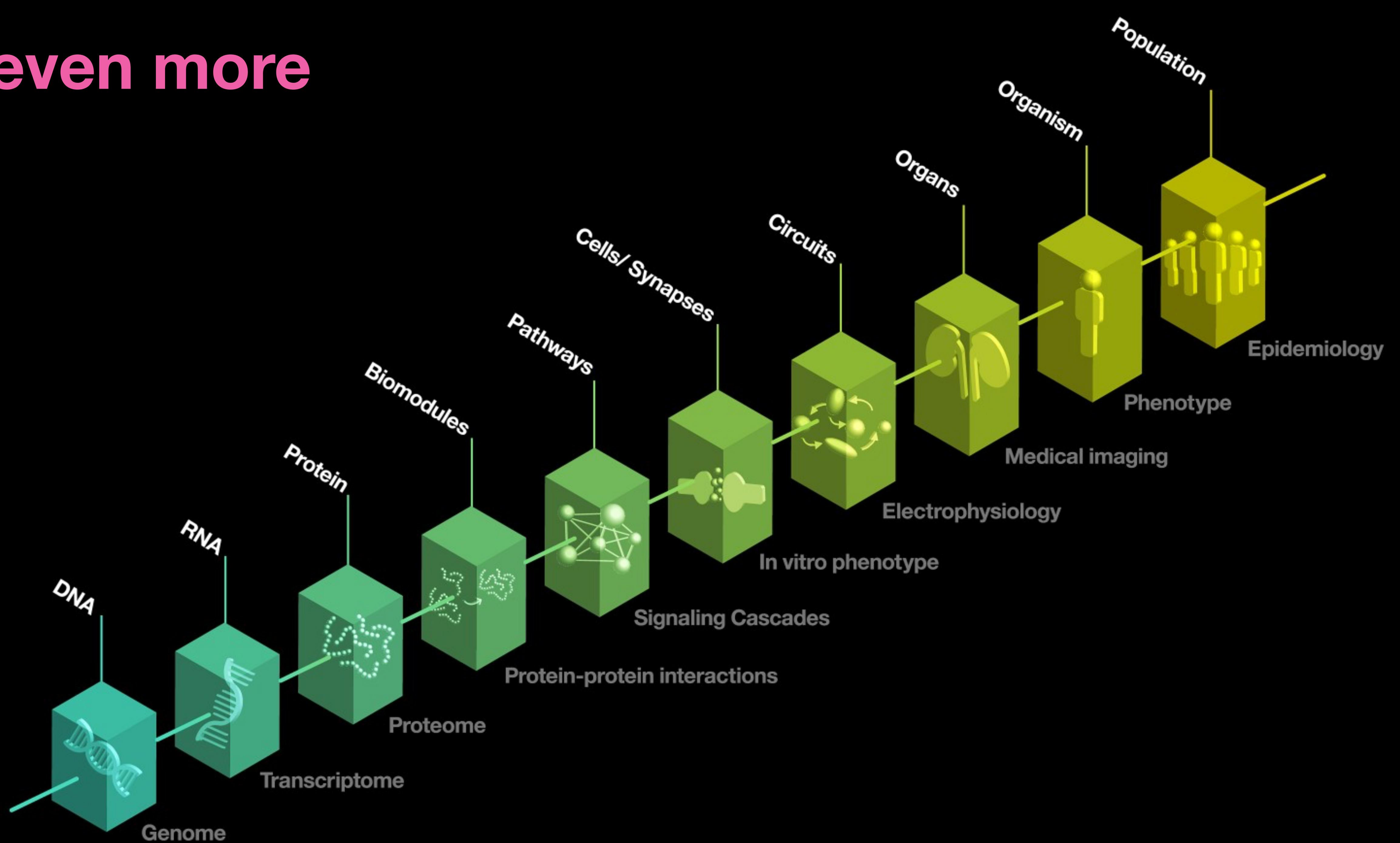
visual: IBM





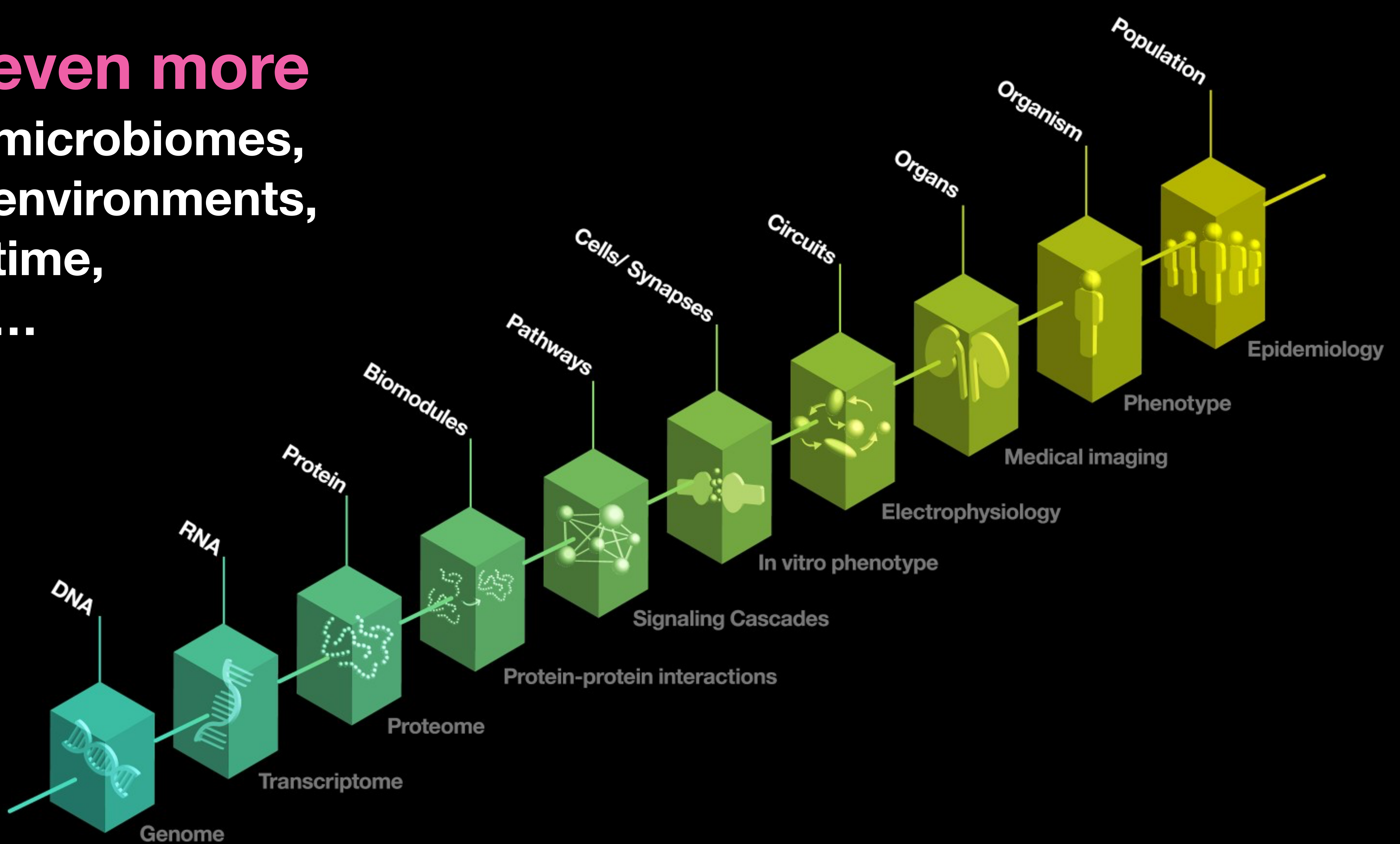


even more



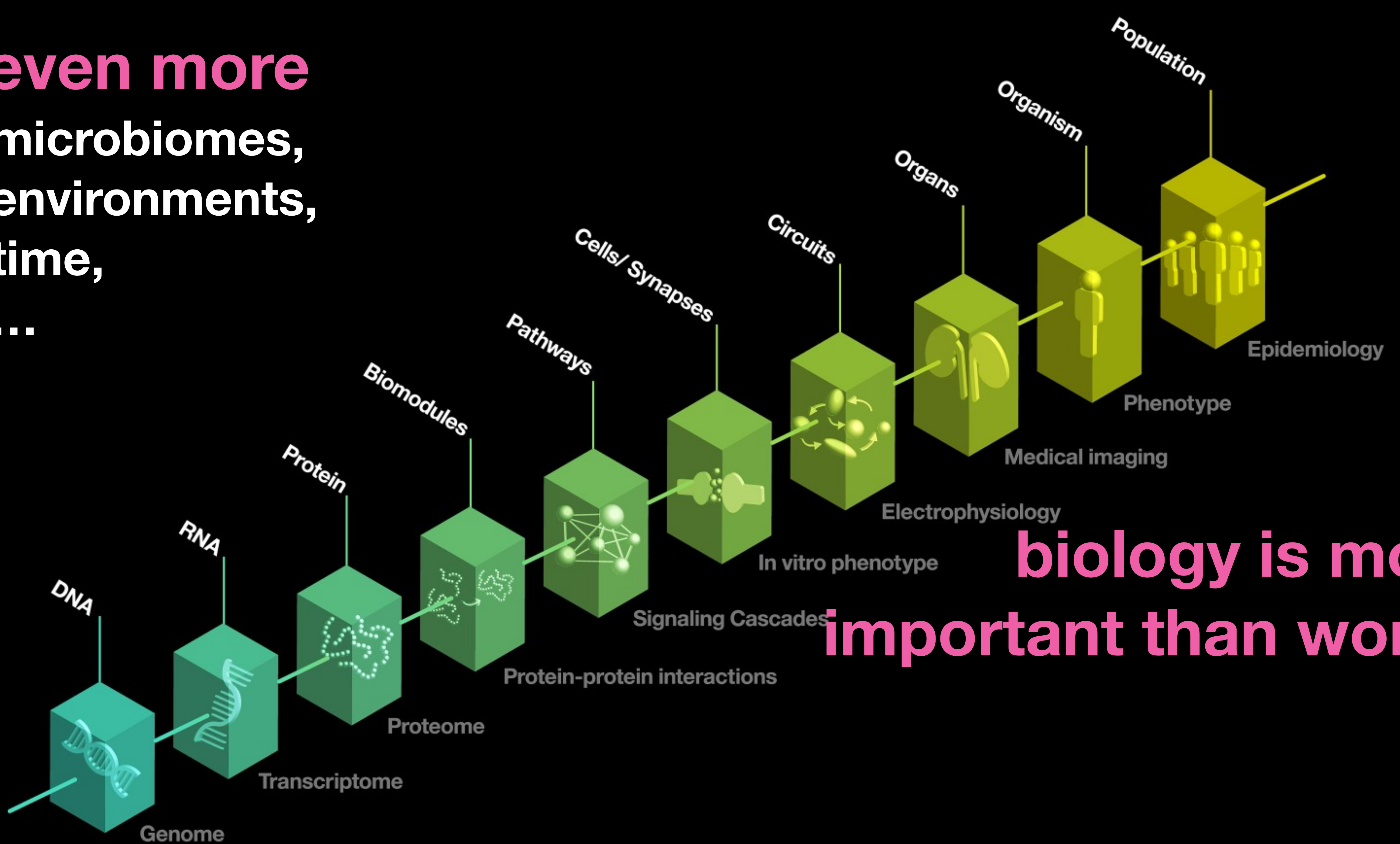
even more
microbiomes,
environments,
time,

...



even more
microbiomes,
environments,
time,

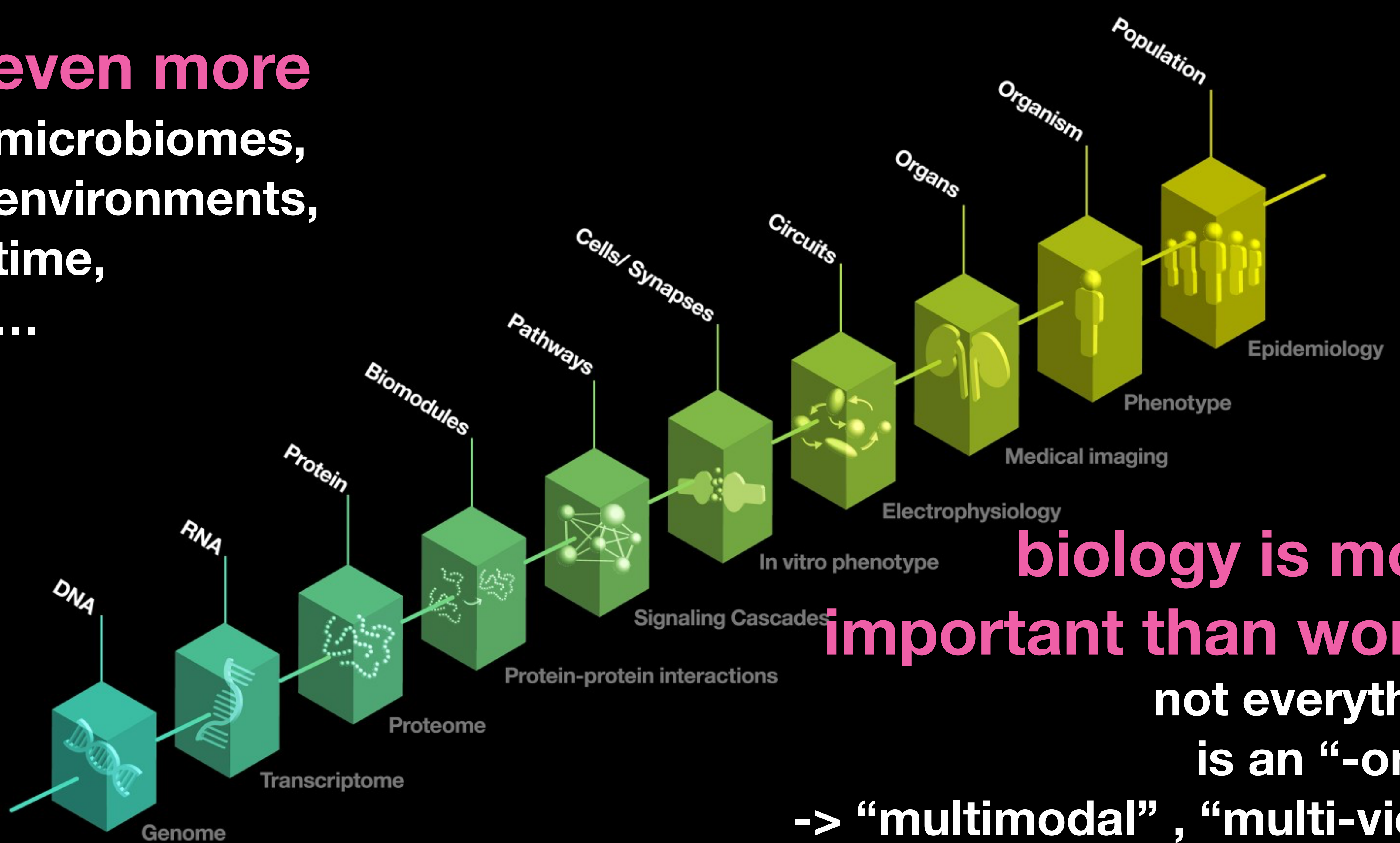
...



biology is more
important than words

even more
microbiomes,
environments,
time,

...



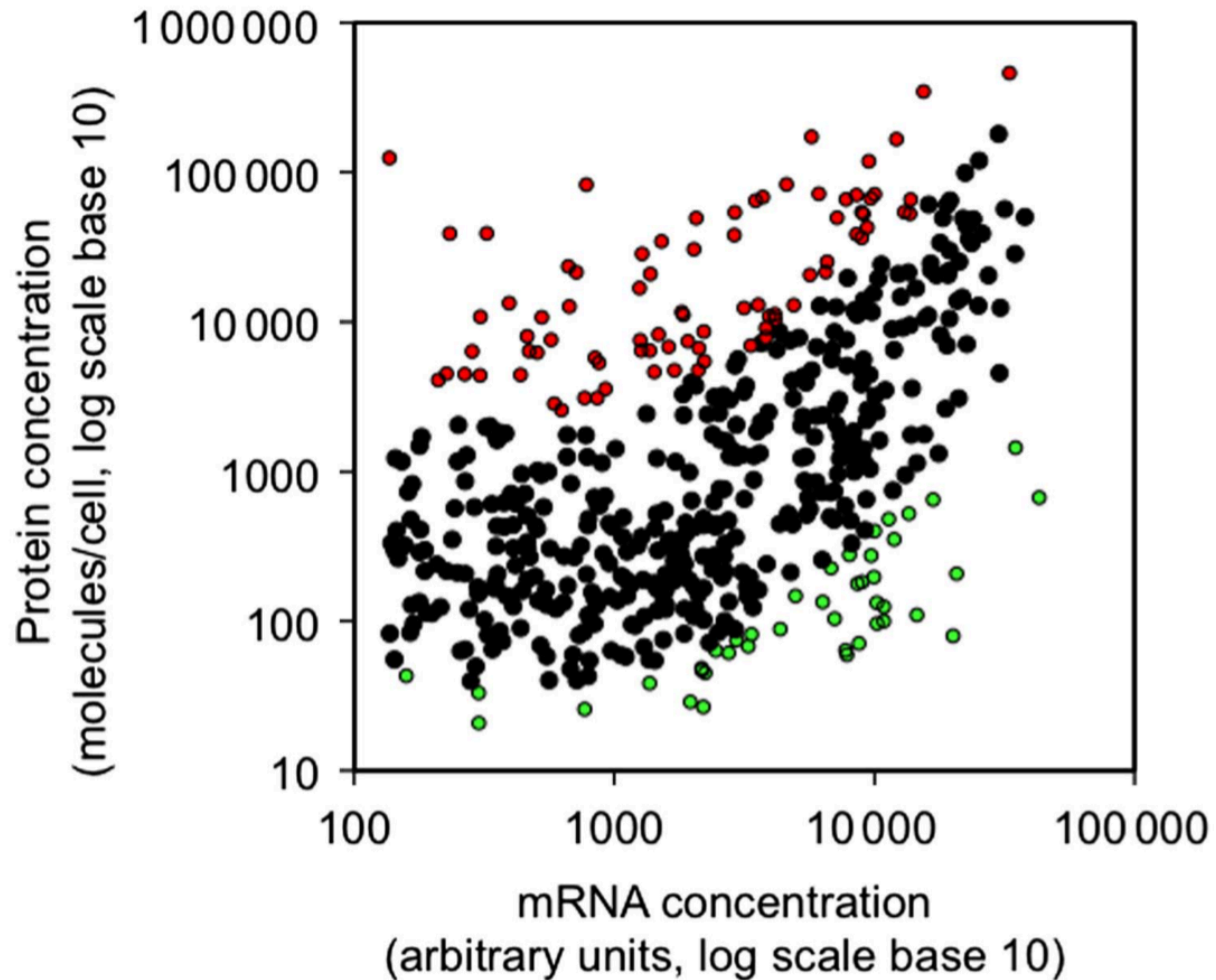
**biology is more
important than words**

**not everything
is an “-ome”**

-> “multimodal” , “multi-view”

many combinations

Be careful!



There are many different strategies.

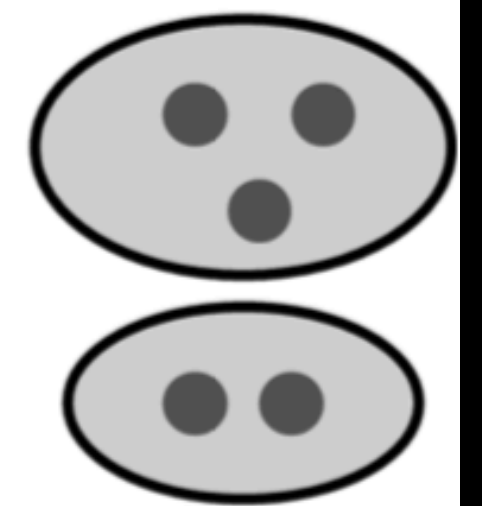
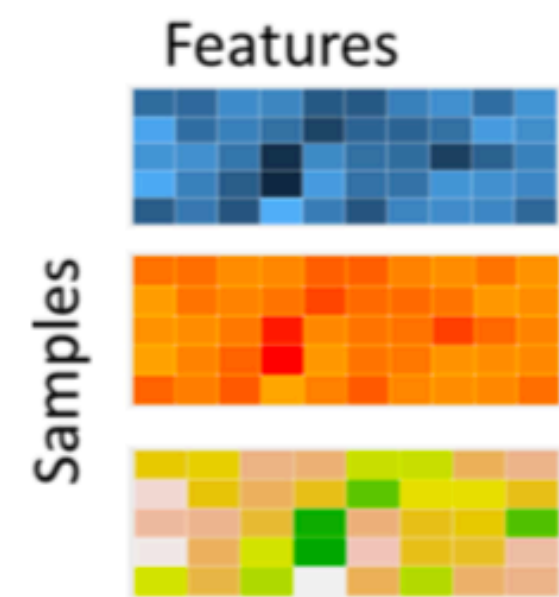


Figure 1. Overview of multi-omics clustering approaches.

Pavlidis et al., 2001
Rappoport and Shamir, 2018
Stuart and Satija, 2019

There are many different strategies.

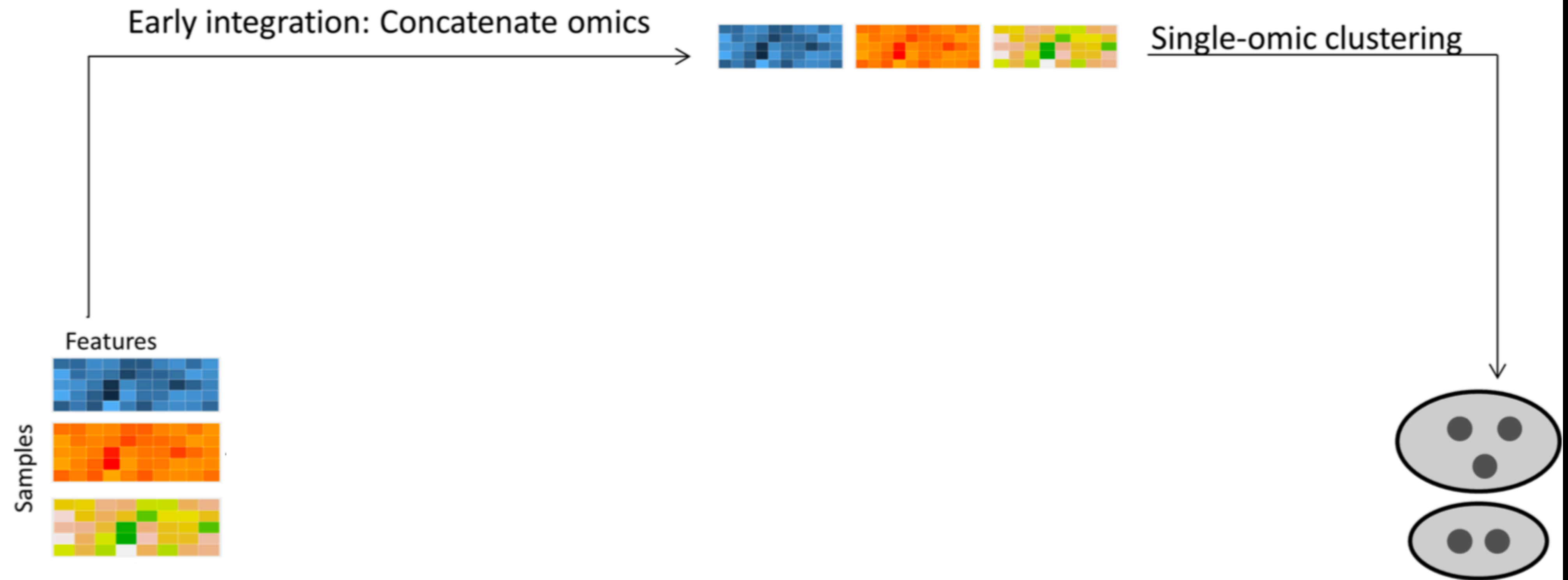


Figure 1. Overview of multi-omics clustering approaches.

Pavlidis et al., 2001
Rappoport and Shamir, 2018
Stuart and Satija, 2019

There are many different strategies.

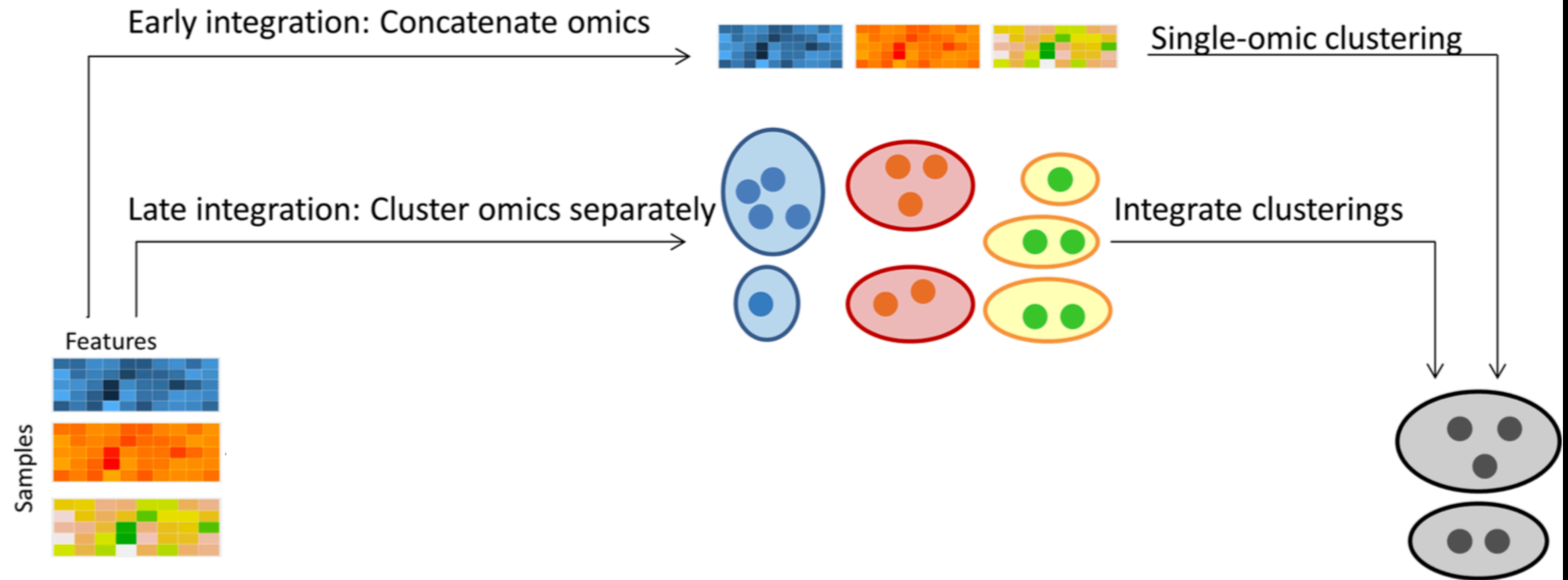
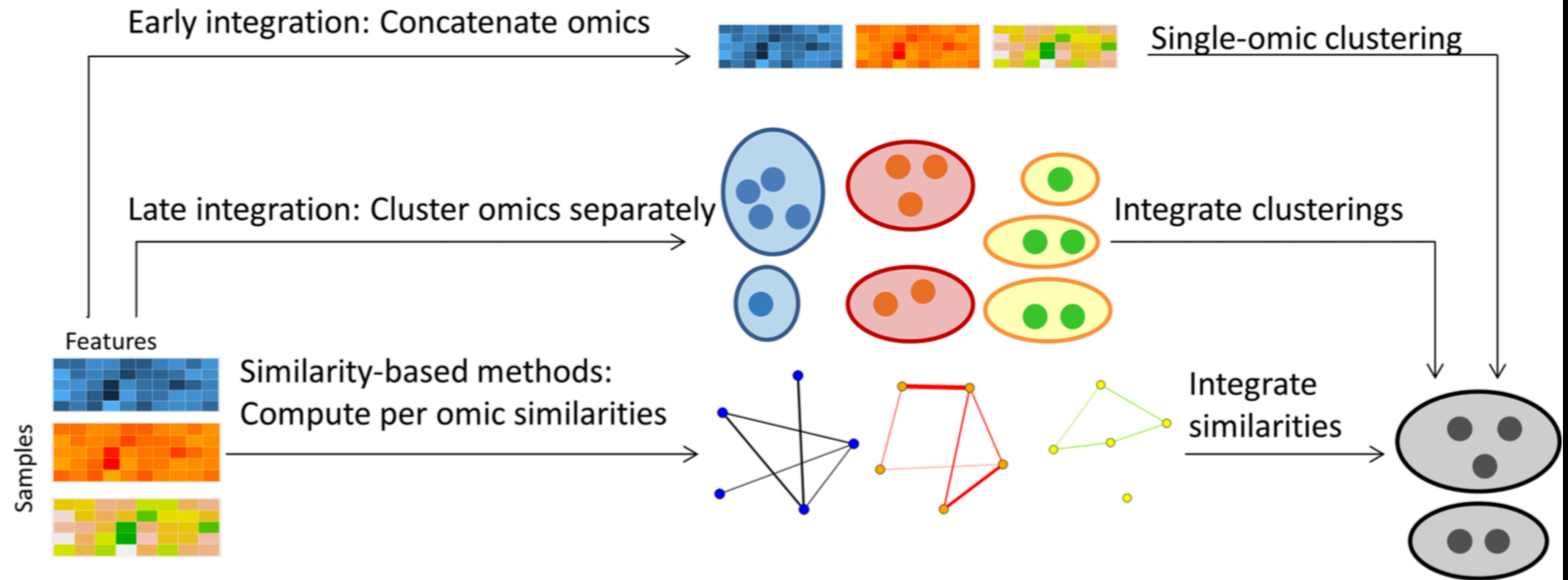


Figure 1. Overview of multi-omics clustering approaches.

Pavlidis et al., 2001
Rappoport and Shamir, 2018
Stuart and Satija, 2019

There are many different strategies.



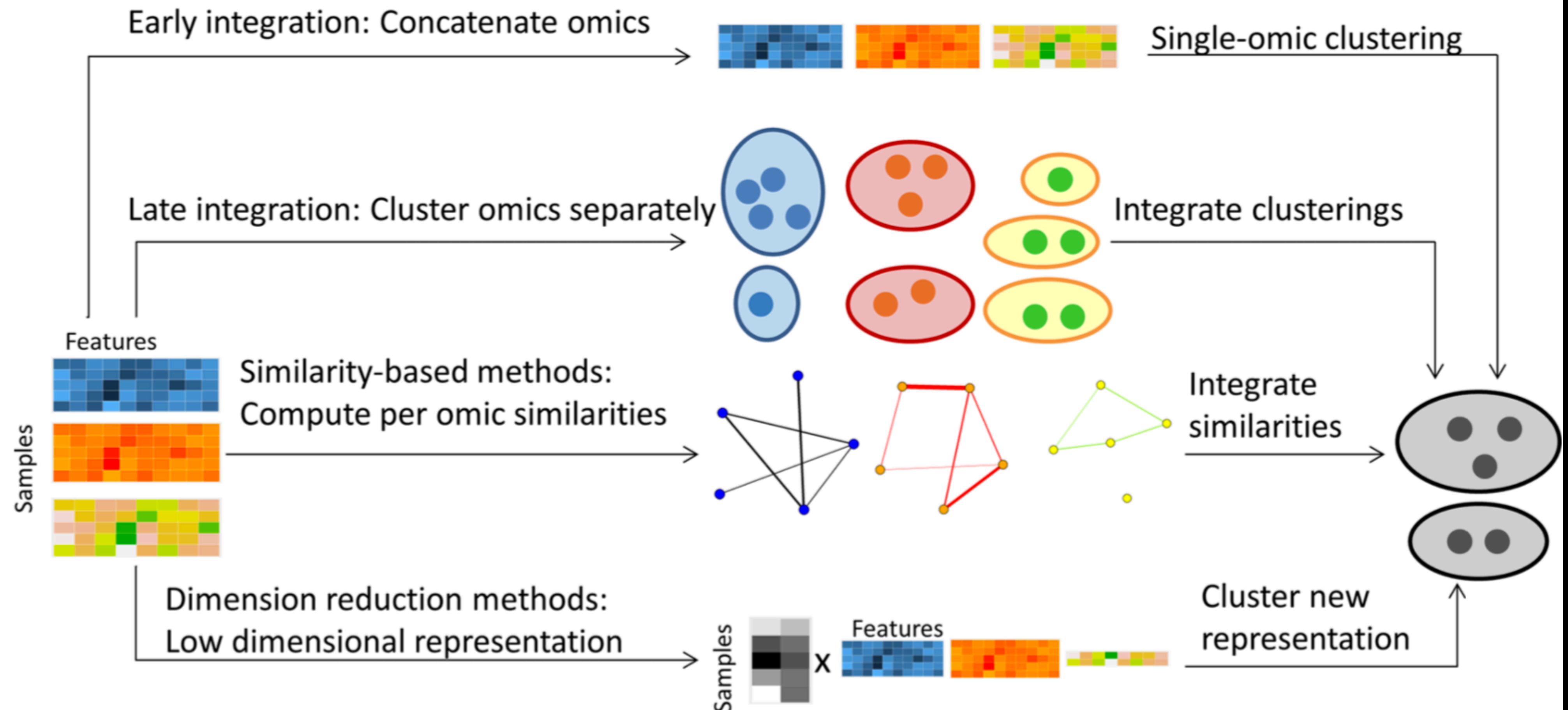
Pavlidis et al., 2001

Rappoport and Shamir, 2018

Stuart and Satija, 2019

Figure 1. Overview of multi-omics clustering approaches.

There are many different strategies.



Pavlidis et al., 2001

Rappoport and Shamir, 2018

Stuart and Satija, 2019

Figure 1. Overview of multi-omics clustering approaches.

Organize data by adhering to good practices.

THE AMERICAN STATISTICIAN

2018, VOL. 72, NO. 1, 2–10

<https://doi.org/10.1080/00031305.2017.1375989>



Taylor & Francis
Taylor & Francis Group

OPEN ACCESS



Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA

ABSTRACT

Spreadsheets are widely used software tools for data entry, storage, analysis, and visualization. Focusing on the data entry and storage aspects, this article offers practical recommendations for organizing spreadsheet data to reduce errors and ease later analyses. The basic principles are: be consistent, write dates like YYYY-MM-DD, do not leave any cells empty, put just one thing in a cell, organize the data as a single rectangle (with subjects as rows and variables as columns, and with a single header row), create a data dictionary, do not include calculations in the raw data files, do not use font color or highlighting as data, choose good names for things, make backups, use data validation to avoid data entry errors, and save the data in plain text files.

ARTICLE HISTORY

Received June 2017

Revised August 2017

KEYWORDS

Data management; Data organization; Microsoft Excel; Spreadsheets

Wrappers, something very useful to learn
for everyone working with data.

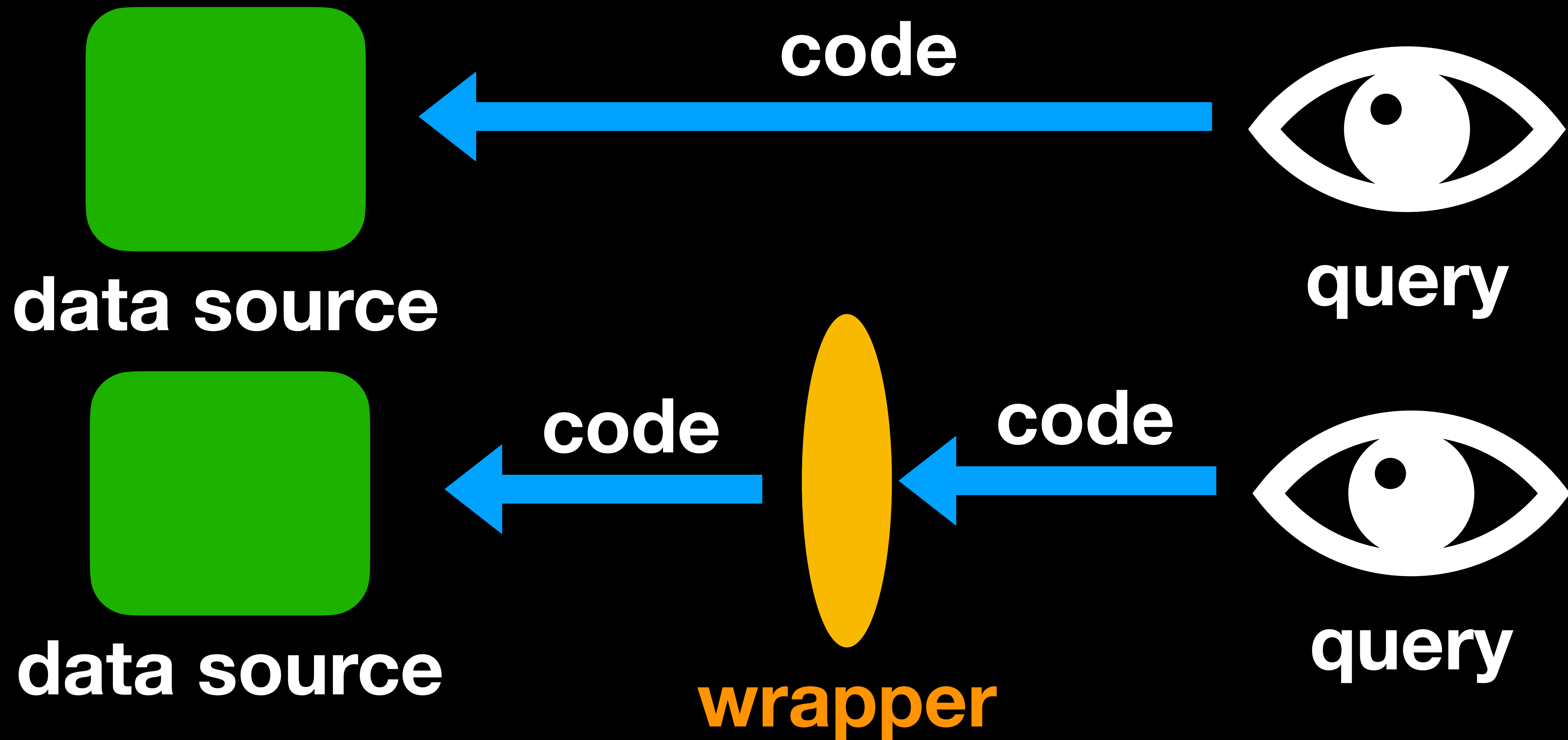


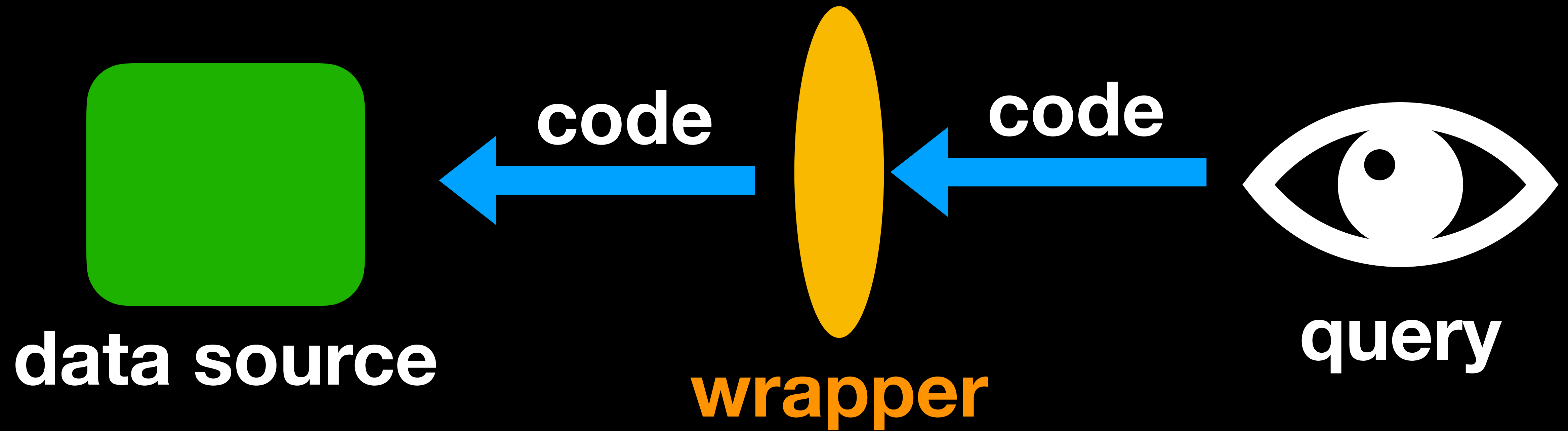
data source

code



query





other / new
differently
formatted



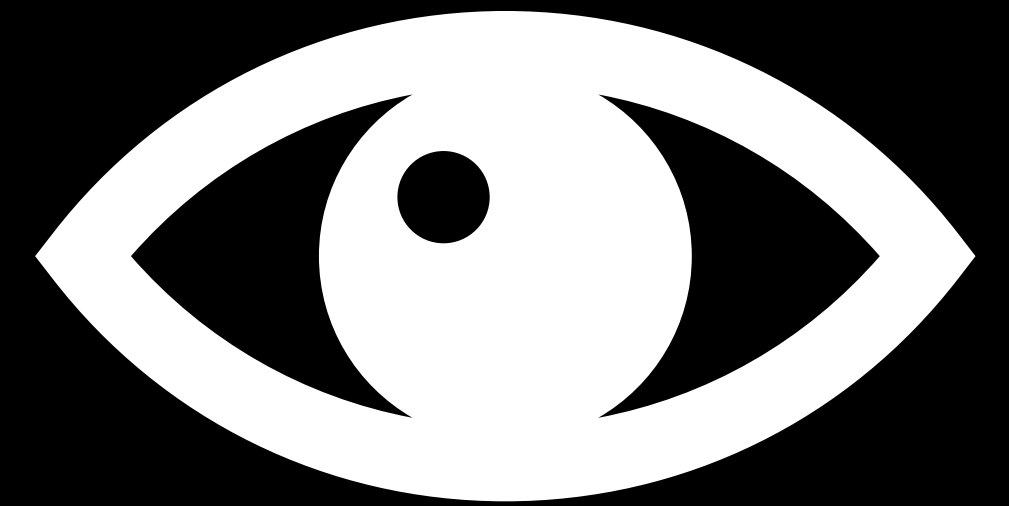
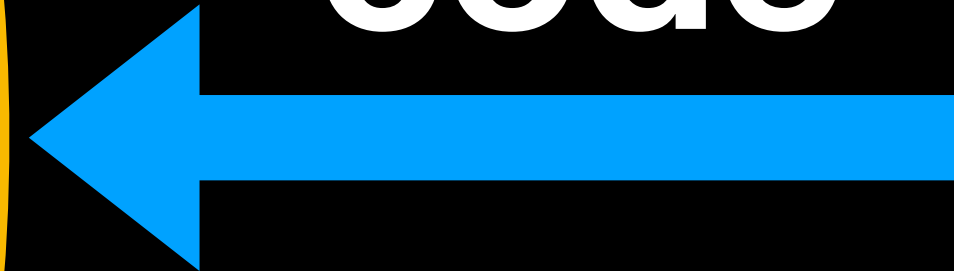
data source

code



wrapper

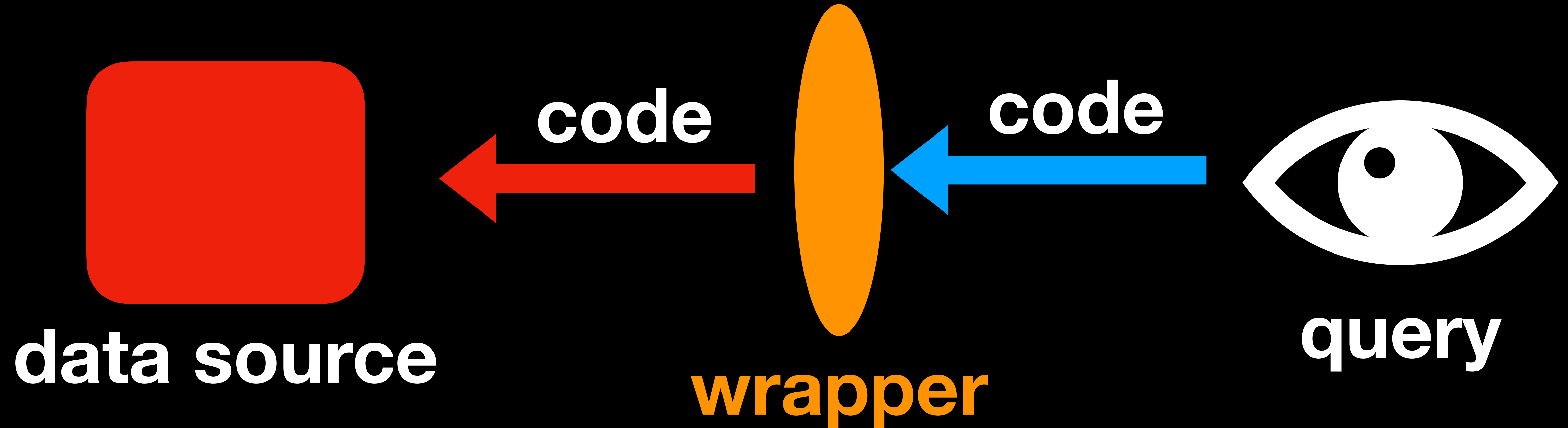
code

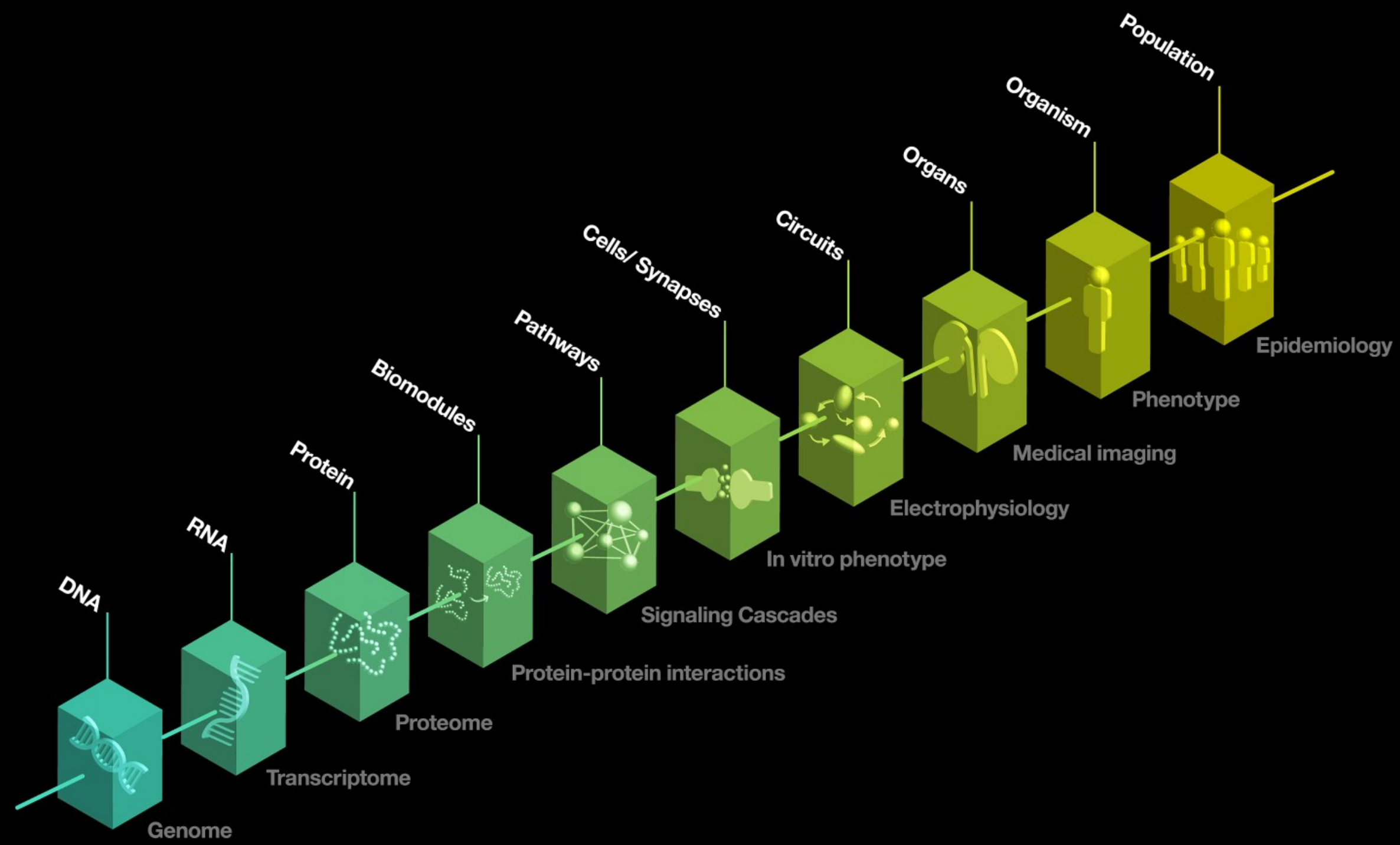


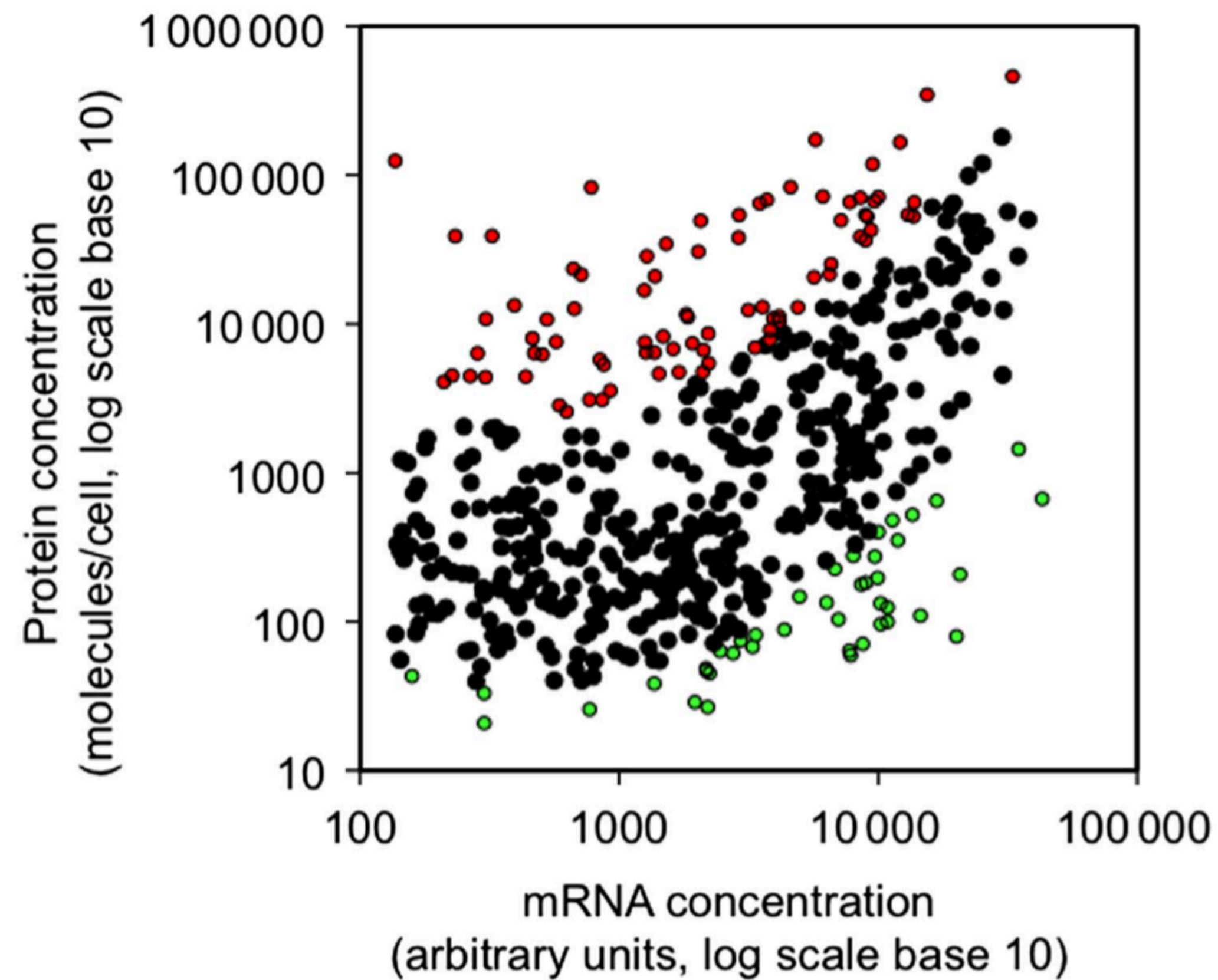
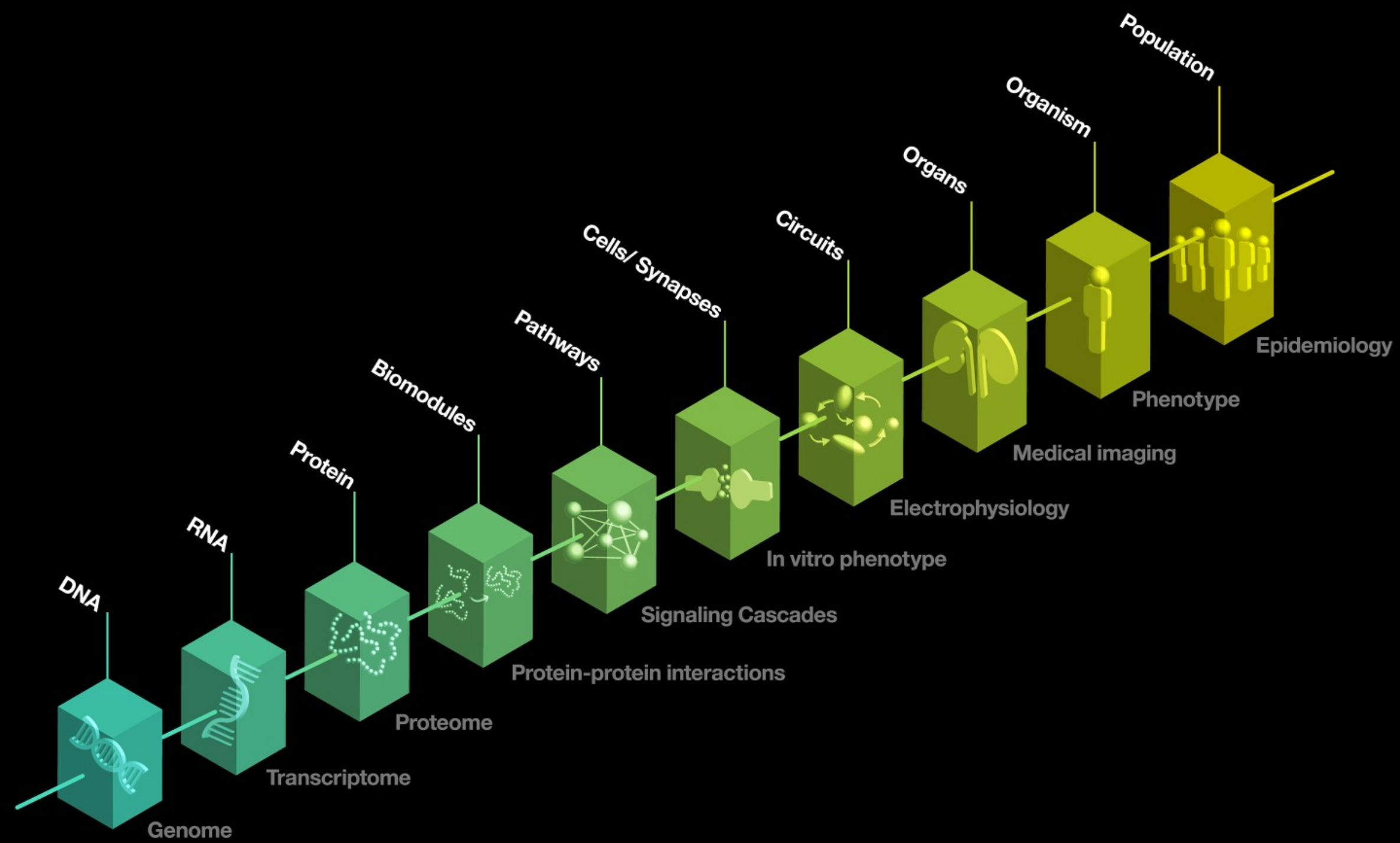
query

other / new
differently
formatted

optimize
speed if
frequently used







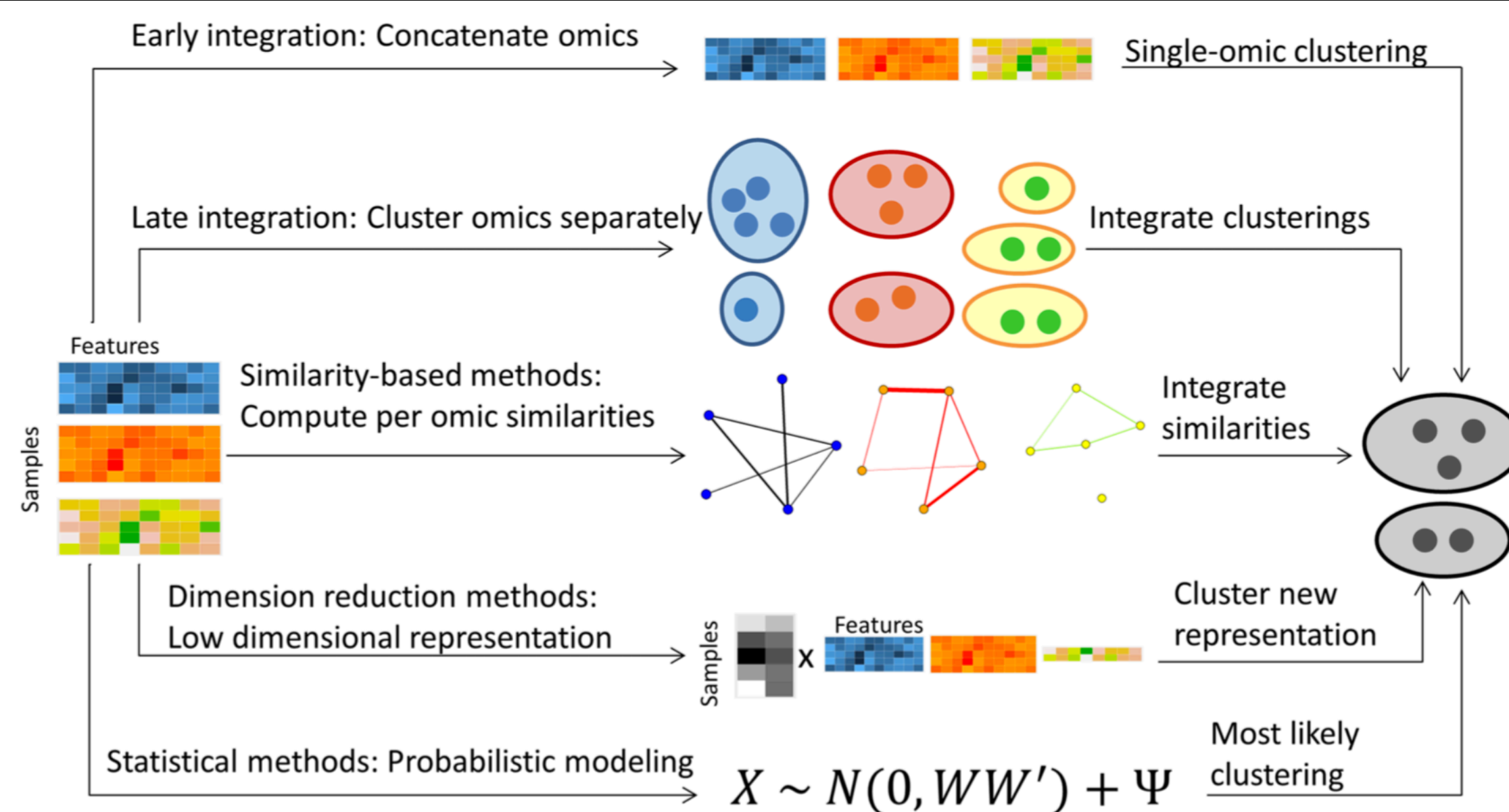
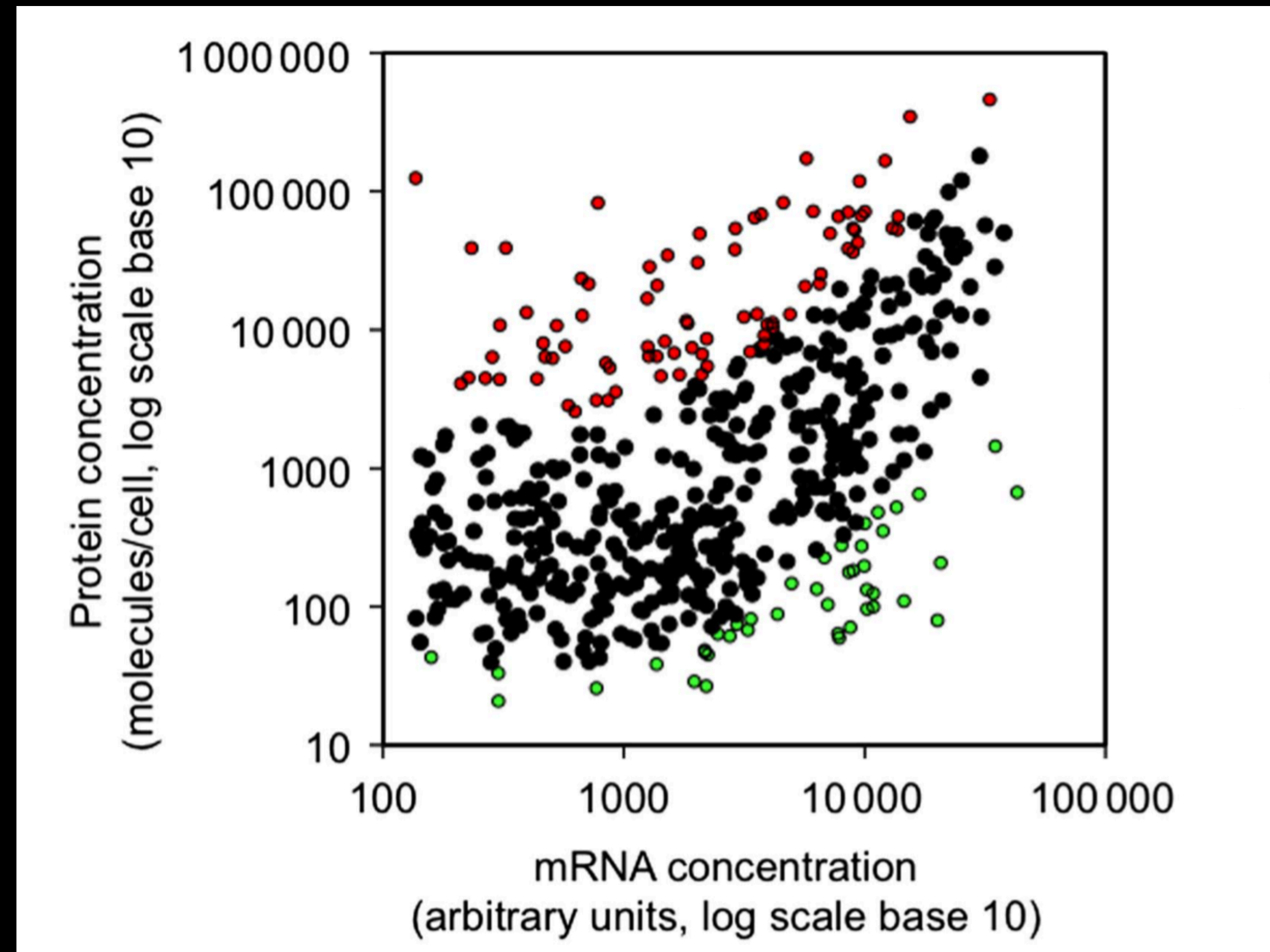
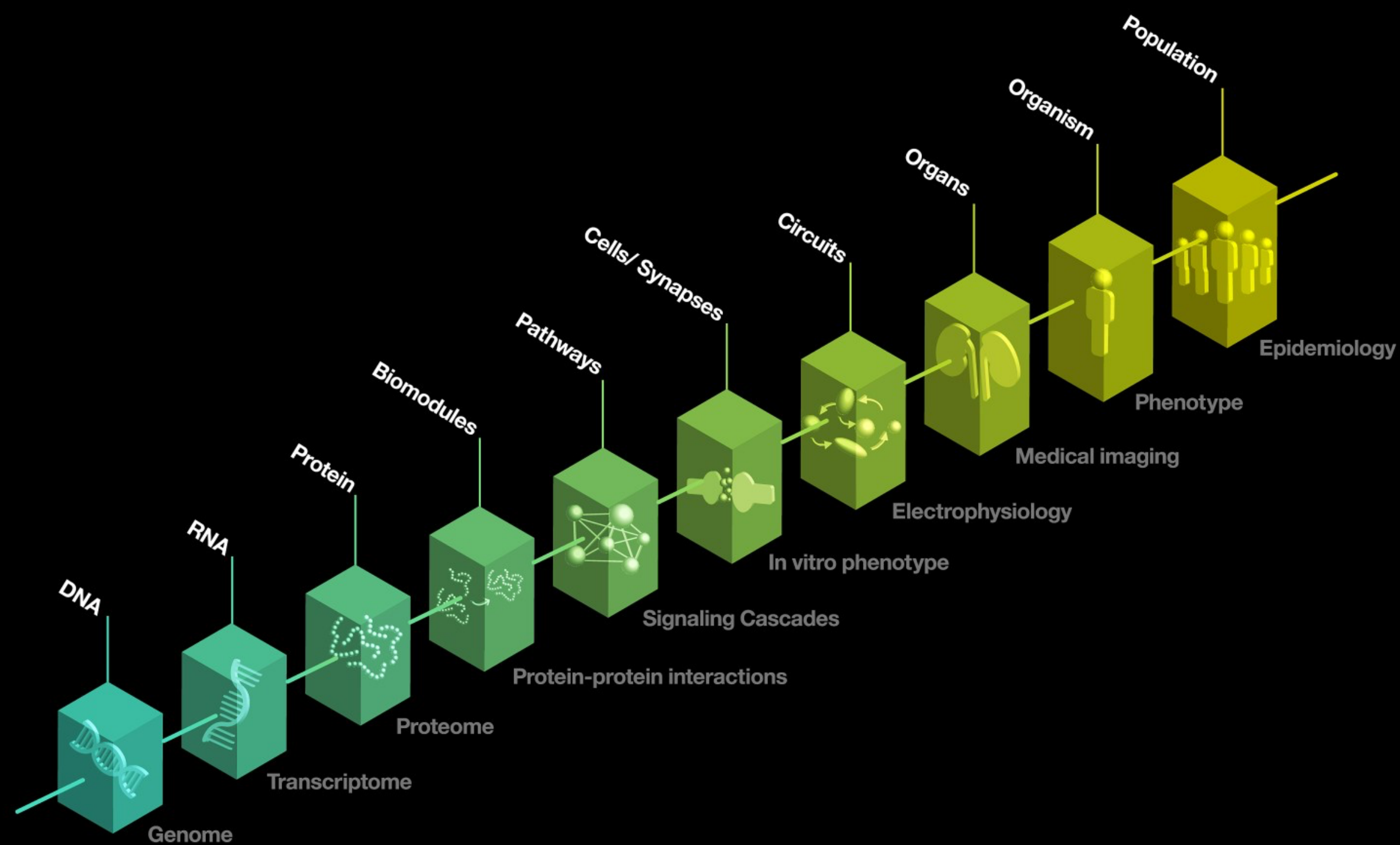
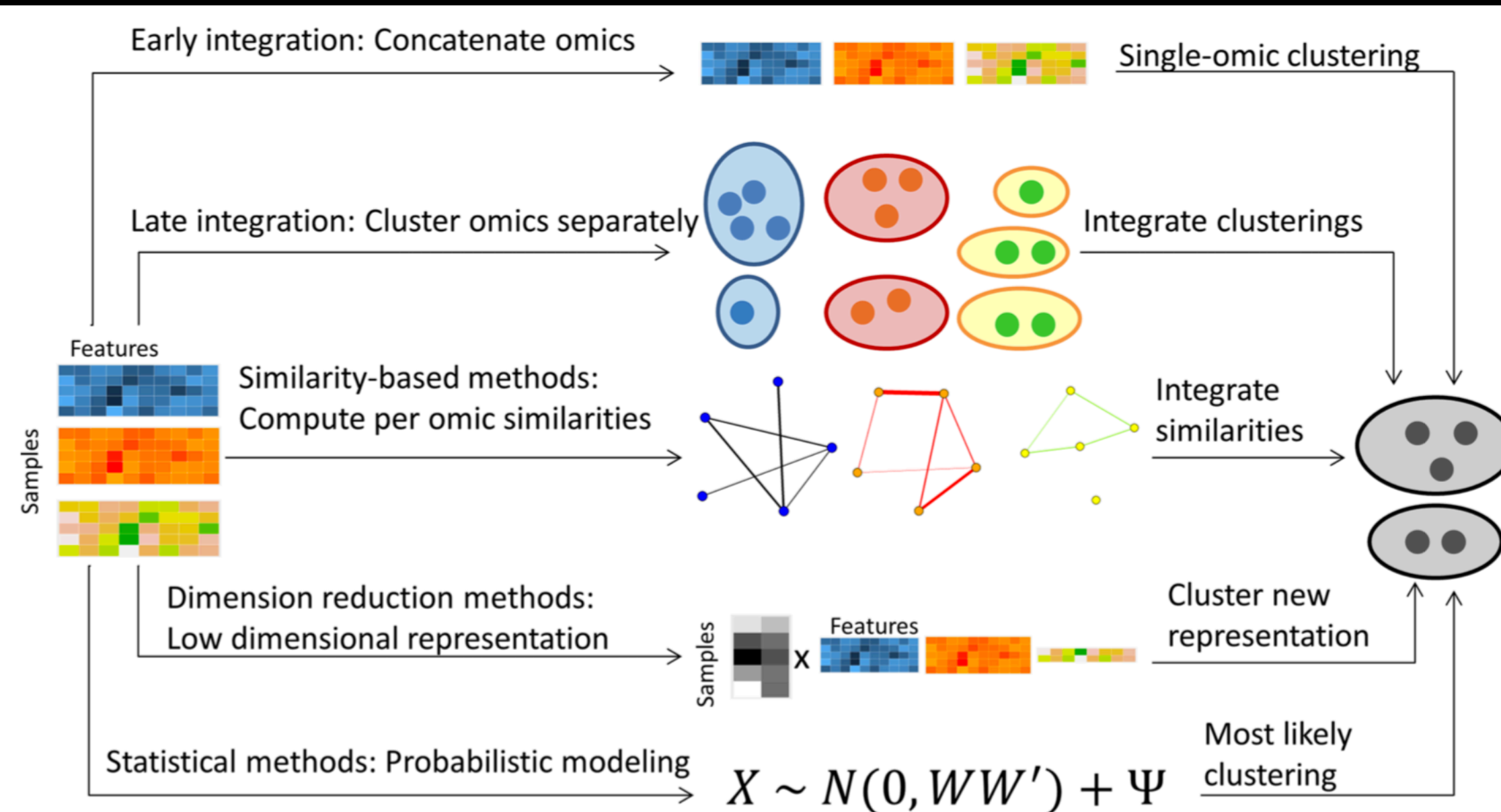
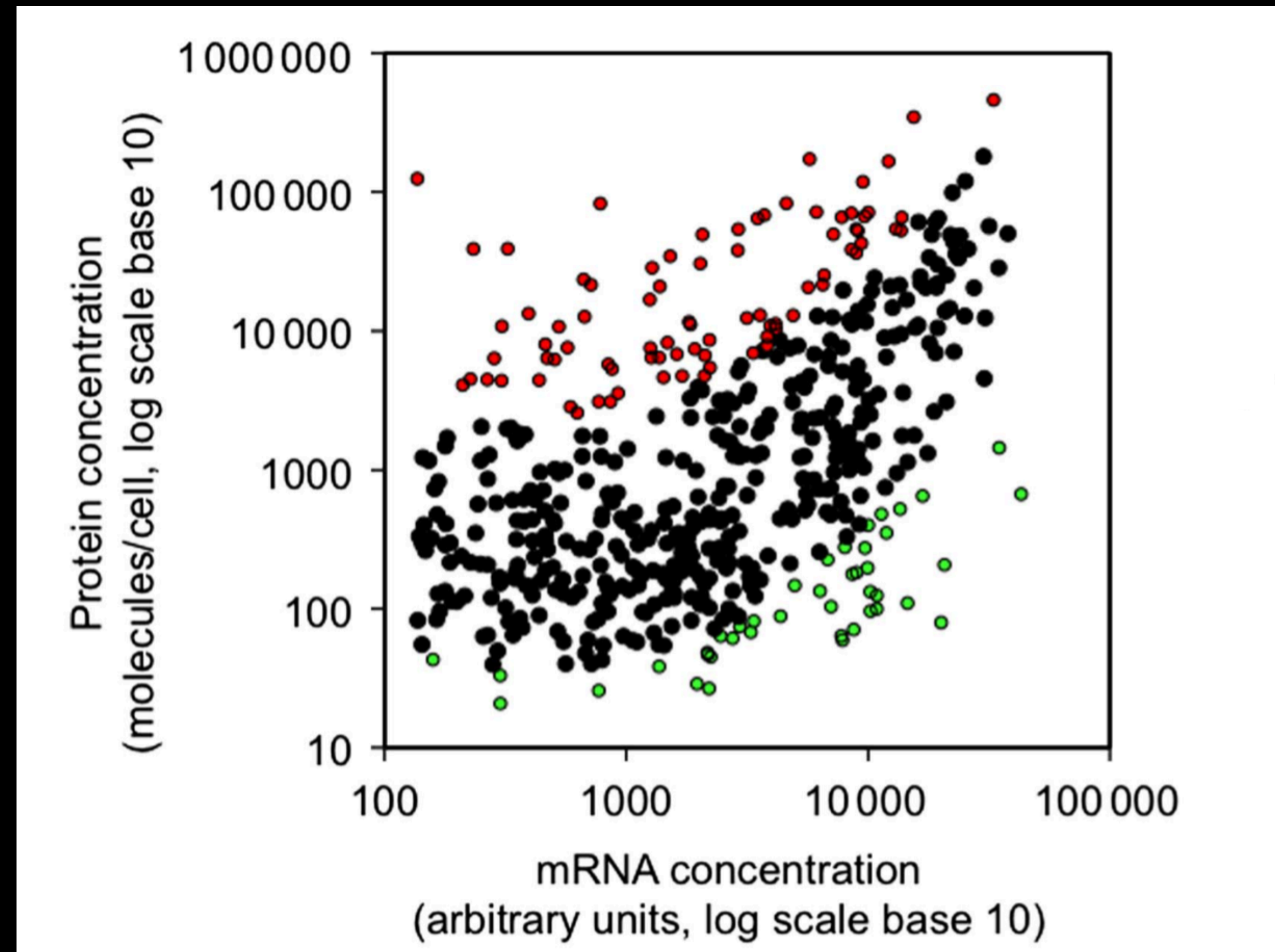
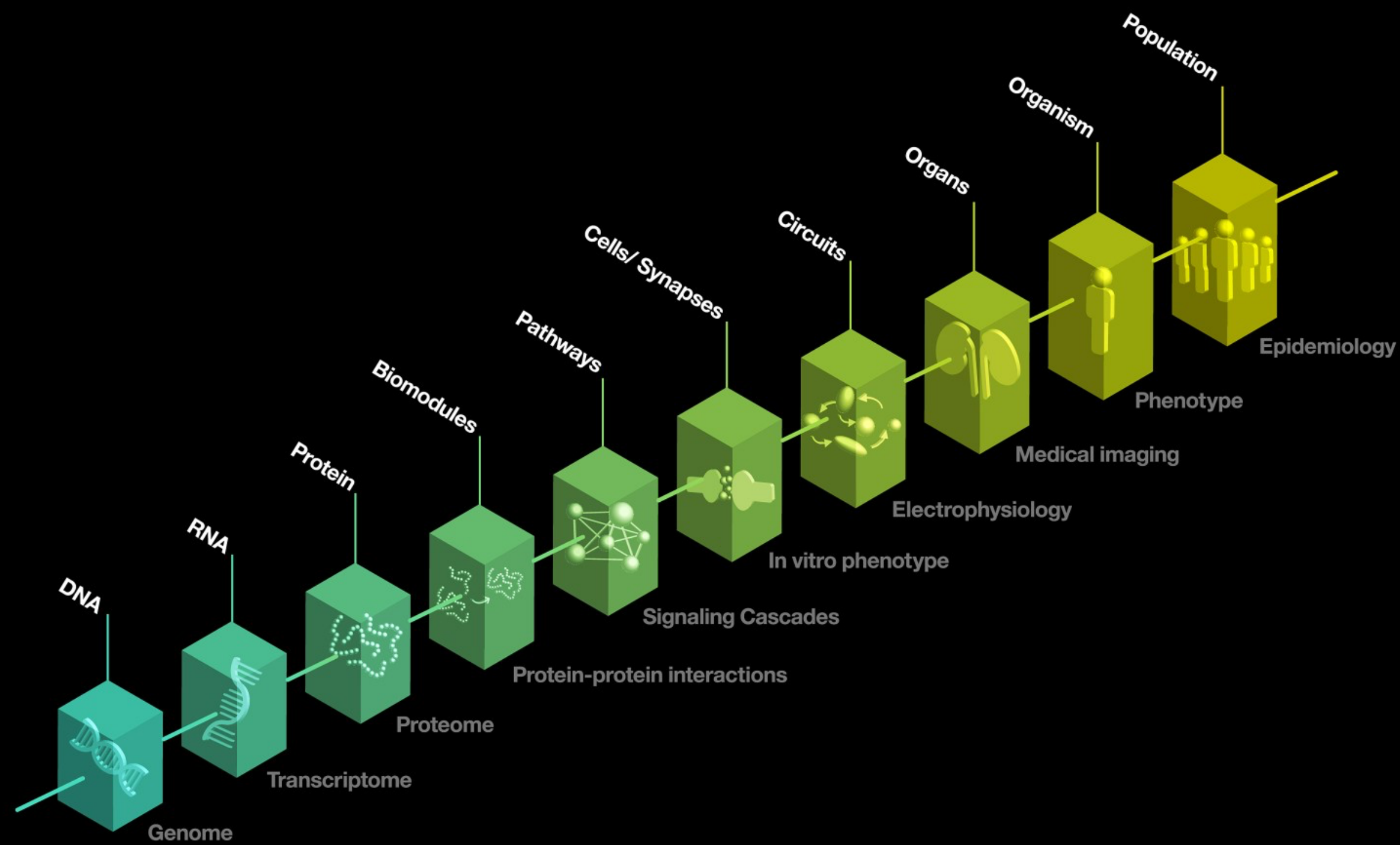


Figure 1. Overview of multi-omics clustering approaches.



← wrapper →

Figure 1. Overview of multi-omics clustering approaches.

group coding exercise experience

30 minutes

30 minutes

3 people per breakout group

30 minutes

3 people per breakout group

free choice of tools

30 minutes

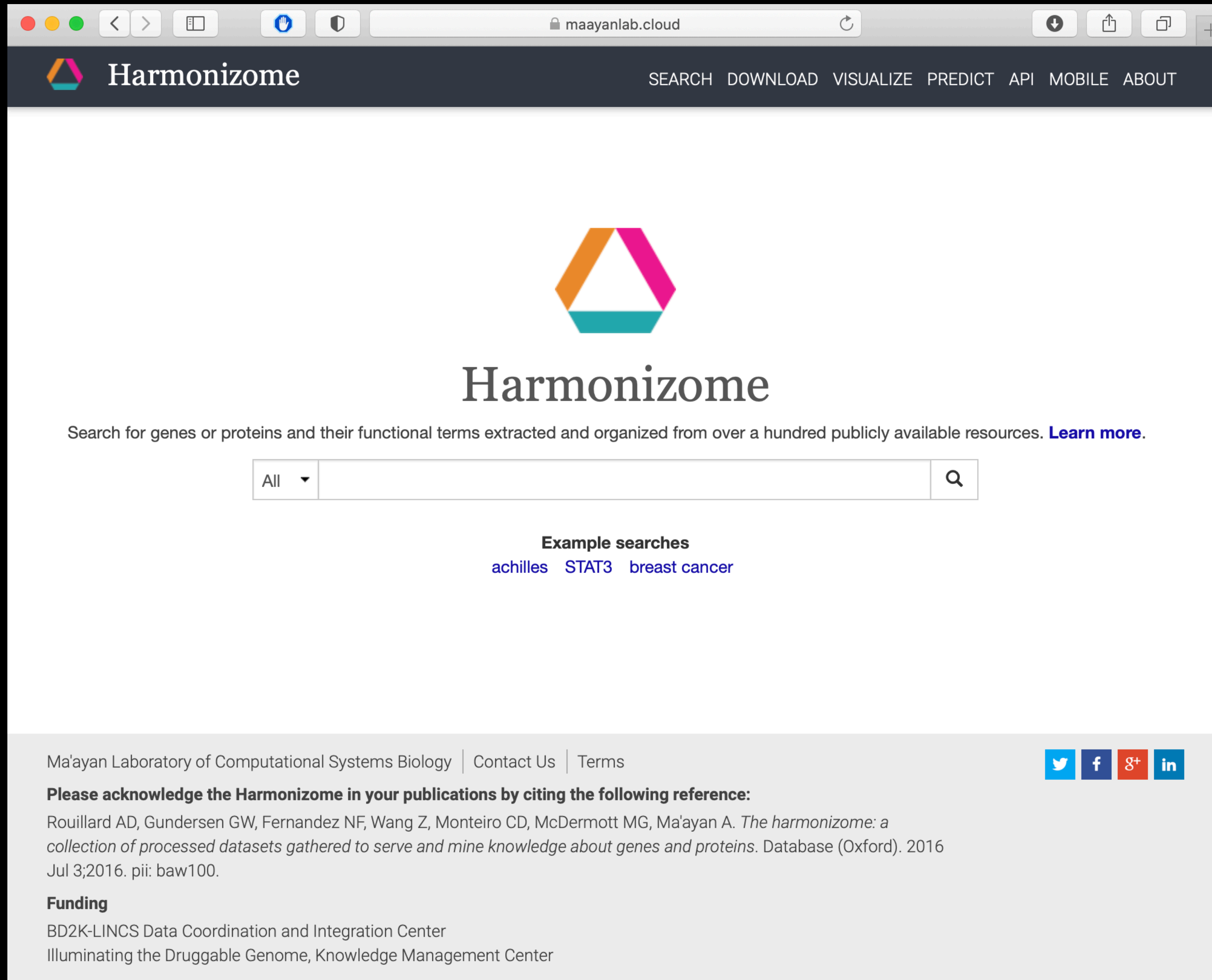
3 people per breakout group

free choice of tools

Are similarities between genes conserved
across scales?

Data Source: <https://maayanlab.cloud/Harmonizome/>

subset: <https://northwestern.box.com/s/dvrxd7ioe6jgm2srrs7rj8mzkqpkhsa> (see handout)



The screenshot shows the Harmonizome website in a web browser. The browser's address bar displays "maayanlab.cloud". The website's header features the Harmonizome logo on the left and a navigation menu with links: SEARCH, DOWNLOAD, VISUALIZE, PREDICT, API, MOBILE, and ABOUT. The main content area has a large Harmonizome logo and the text "Harmonizome". Below this, a search bar is present with a dropdown menu set to "All" and a search icon. A section titled "Example searches" lists "achilles", "STAT3", and "breast cancer". The footer contains contact information for the Ma'ayan Laboratory of Computational Systems Biology, a link to "Contact Us", and "Terms". It also includes a request to acknowledge the Harmonizome in publications with a reference to Rouillard et al. (2016). Social media icons for Twitter, Facebook, Google+, and LinkedIn are displayed. The "Funding" section lists the BD2K-LINCS Data Coordination and Integration Center and the Illuminating the Druggable Genome, Knowledge Management Center.

Ma'ayan Laboratory of Computational Systems Biology | [Contact Us](#) | [Terms](#)

Please acknowledge the Harmonizome in your publications by citing the following reference:
Rouillard AD, Gundersen GW, Fernandez NF, Wang Z, Monteiro CD, McDermott MG, Ma'ayan A. *The harmonizome: a collection of processed datasets gathered to serve and mine knowledge about genes and proteins*. Database (Oxford). 2016 Jul 3;2016. pii: baw100.

Funding
BD2K-LINCS Data Coordination and Integration Center
Illuminating the Druggable Genome, Knowledge Management Center

Provides precomputed similarities between genes.

Are similarities between genes conserved across scales?

Note: while you might answer this within 30 minutes, identifying possible challenges is even more valuable.

Suggested: a) take two datasets that intrigue you – possibly small to limit download time b) work / talk with your fellow team members

Data Source: <https://maayanlab.cloud/Harmonizome/>

subset: <https://northwestern.box.com/s/dvrxd7ioe6jgm2srrs7rj8mzkqpkhsa> (see handout)

Template code: https://github.com/tstoeger/course_multi_omics/blob/main/how_well_do_different_views_on_biology_correlate.ipynb

Ask anonymously on: <https://padlet.com/thomasstoeger/a4mtvpym671dwgnl>