

GraphDB's & Applications

Denis Vrdoljak | Gunnar Kleemann

UC Berkeley School of information
Berkeley Data Science Group, LLC



Presentation Road Map

- Intro
- Background
- Examples
- Our Work
- Graph Databases



Intro

Background

Examples

Our Work

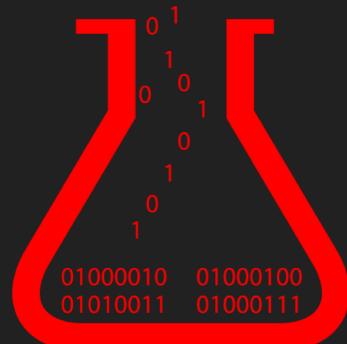
Graph Databases



About Us: BDSG

Berkeley Data Science Group

Founded by UC Berkeley Data Science instructors and alumni with the goal of bringing Berkeley data science projects to market and commercializing Berkeley Data Science research.



About Us: the Speakers

Denis Vrdoljak

denis@bds.group

dvradolja@cisco.com

Gunnar Kleemann, PhD
gunnar@bds.group



Why Graphs?

Why Graph Databases?



GraphDB's Optimized for Relationships

- Graph Databases store data in tables/rows/columns, just like a traditional RDBMS
- First Class Citizen is a relationship, not an entity
- GraphDB's are optimized for graph traversals
- This also makes them slow at data retrieval
- But, they're a LOT faster at traversing the nodes of a graph!



Intro

Background

Examples

Our Work

Graph Databases



What is Graph Theory?

Dates back to 1736: Seven Bridges of Königsberg, by Leonhard Euler

Laid down the original groundwork for what became Graph Theory

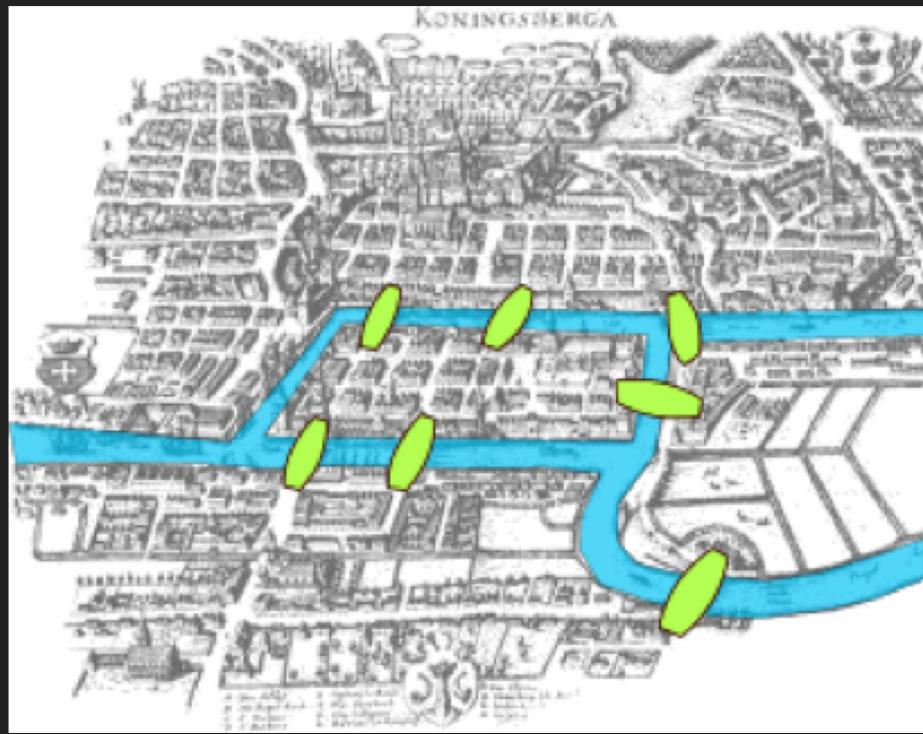
Also evolved into modern day Network Analysis (or Network Graph Analysis) and Social Network Analysis (SNA)



The Seven Bridges of Konigsberg

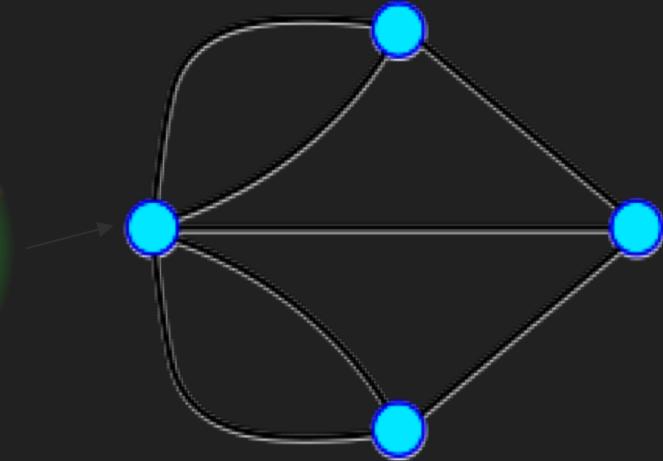
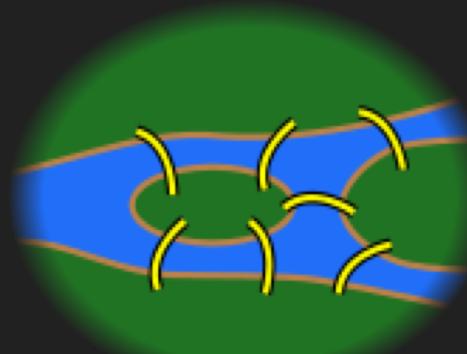
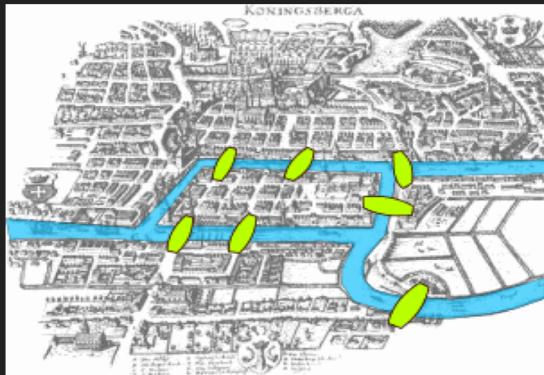
The two large islands in the city of Konigsberg were connected by 7 bridges.

Problem Statement: Can you visit each part of the town, using each bridge only once?



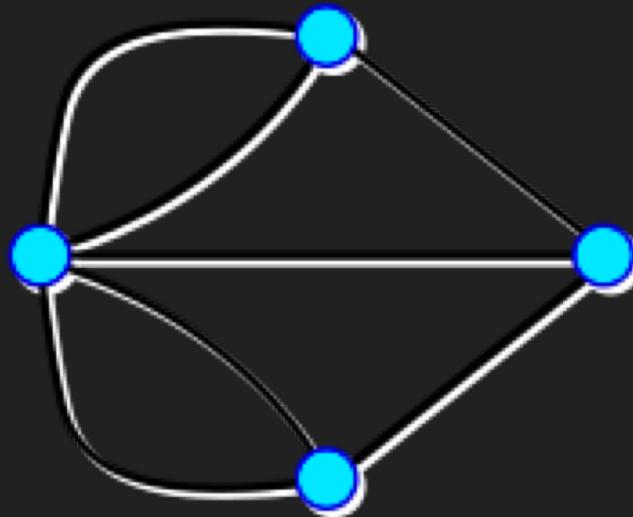
The Seven Bridges of Konigsberg

Leonhardt Euler came up with a new way of thinking about the problem, and in turn became the father of modern Graph Theory



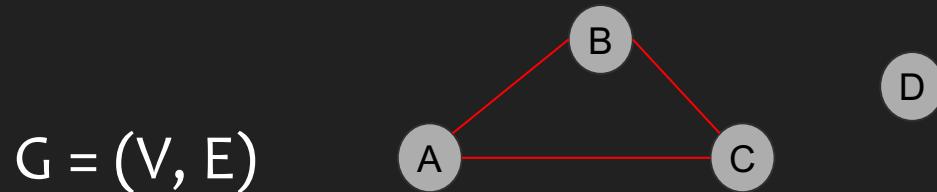
The Seven Bridges of Konigsberg

Euler's Final Answer: it is impossible. His analysis technique laid the foundations of graph theory.



Terminology and Notation

A **graph** G is a set of vertices V that are connected by a set of edges E .



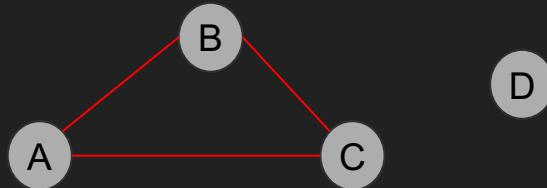
$$V = \{A, B, C, D\}$$

$$E = \{AB, AC, BC\}$$

A vertex A's “children” is the set of vertices that can be reached from vertex A by traversing an edge connected to A.



Terminology and Notation



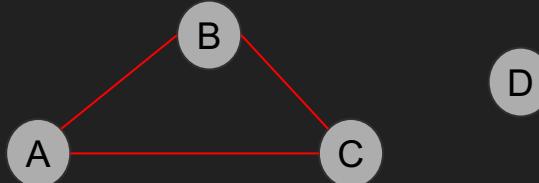
A **cycle** is a group of at least 3 vertices connected in a closed chain (ie, you can start and end a traversal of the vertices at the same vertex).

Vertices A, B, and C form a cycle above.



Terminology and Notation

- Edges of a graph may have **weights** that can be used to model various phenomena
- Edges may also be **directed**, pointing in a particular direction to indicate how the edge may be traversed
 - A graph that does not have any directed edges is **undirected**



Graph Searching Algorithms

Breadth First Search: starting at a vertex N, add all nodes connected to N to a queue, and then continuously breadth first search the items in the queue, in order (adding new vertices to the queue in recursive calls).

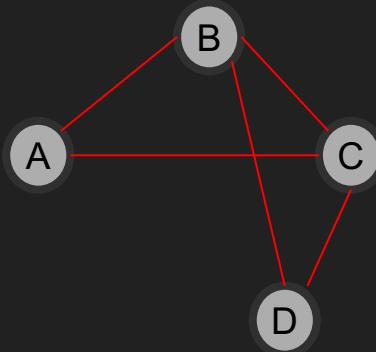
Depth First Search: starting at a vertex N, recursively depth first search all vertices connected to N that have not been searched yet. *Alternatively, depth first search can be performed with the same exact code as breadth first search, replacing the queue with a stack.*



Representations in Code

There are two primary methods of representing graphs in code:

- 1) **Adjacency matrix**: each row / column represents a vertex, and the value in row i and column j represents the weight of the edge between vertices i and j .



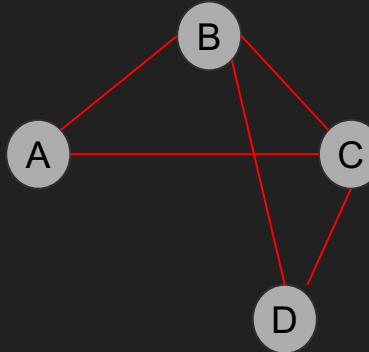
	A	B	C	D
A	0	1	1	0
B	1	0	1	1
C	1	1	0	1
D	0	1	1	0



Representations in Code

There are two primary methods of representing graphs in code:

- 2) **Adjacency list**: A list (usually a linked list) is stored for each vertex, representing its connections.



A: [B, C]

B: [C, D]

C: [A, B, D]

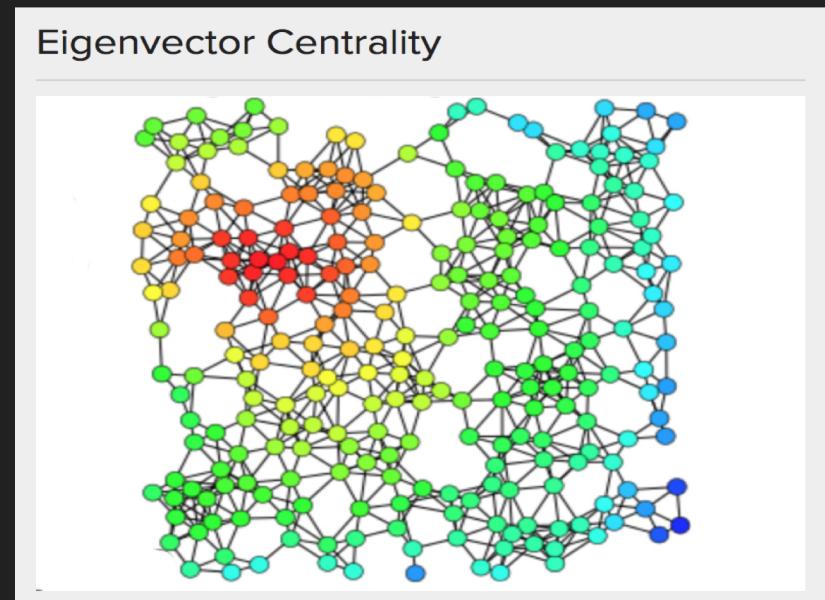
D: [B, C]



Node Level Indices

- Centrality Measures:
 - Degree
 - Betweenness
 - Closeness
 - Eigenvector

Google's PageRank is the quintessential example of Eigenvector Centrality



Picture: Wikipedia

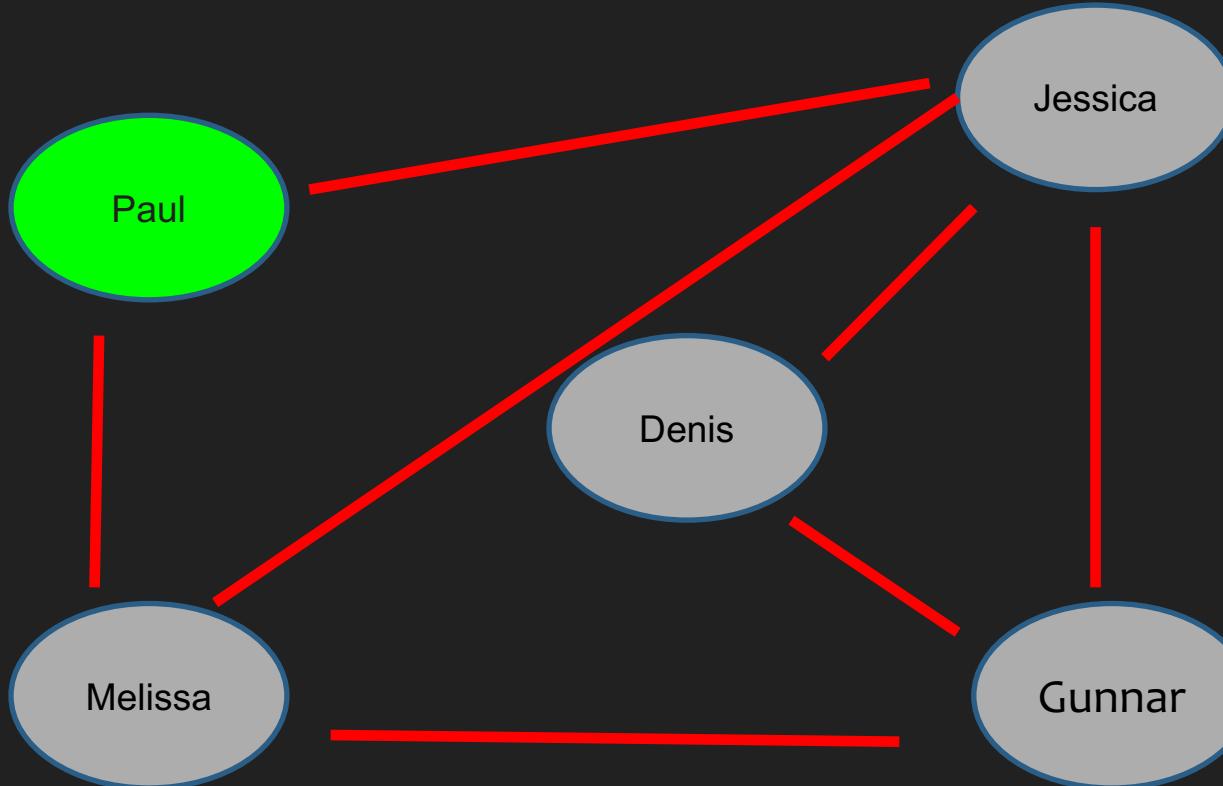


Graph Level Metrics

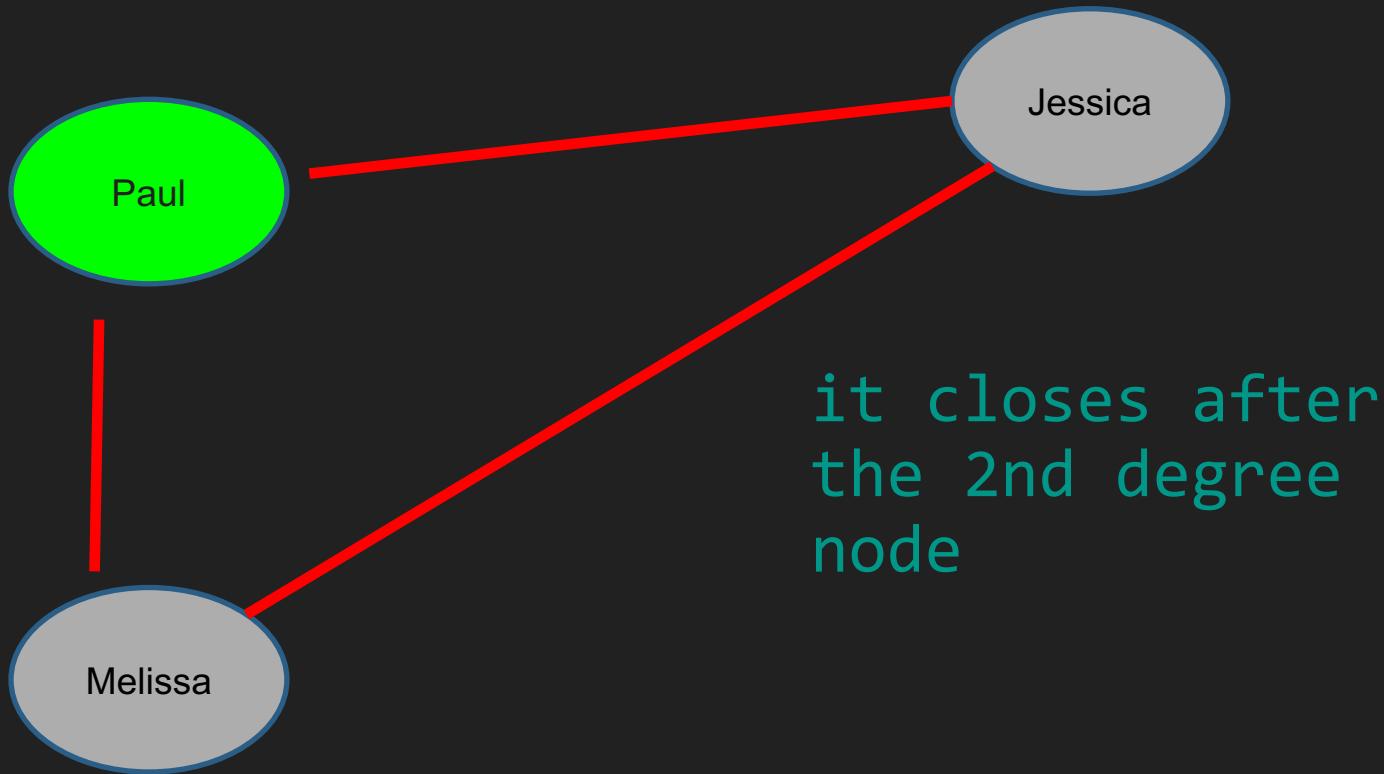
- Averages of Node Level Indices
- Distance
- Community-/Cluster-Based



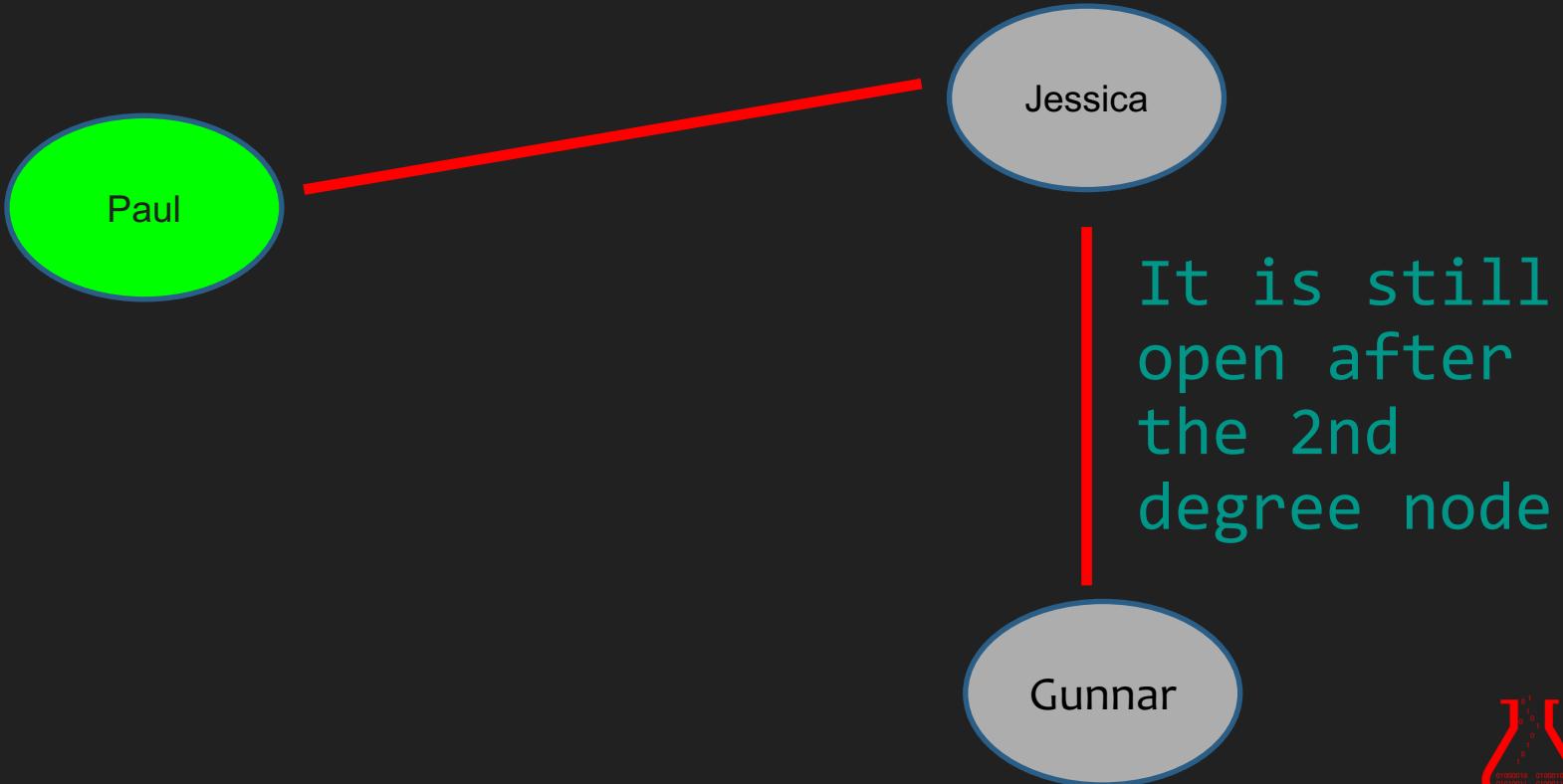
Consider Paul in a simple social network



This relationship forms a triangle



This relationship forms a wedge

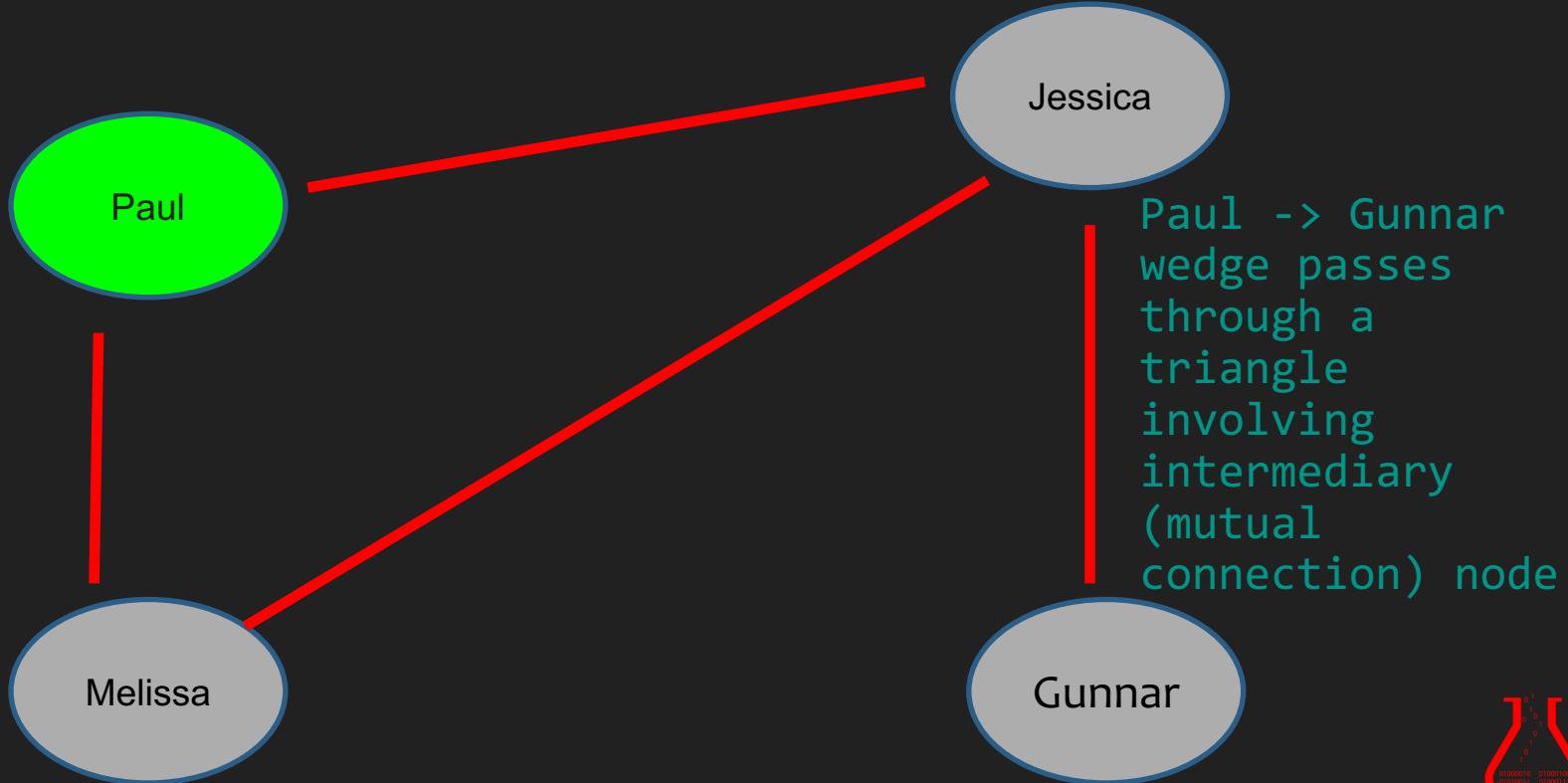


Strength of Wedge

- Number of Triangles that a Wedge pass through, where the intermediate nodes (mutual connection) is part of the triangle.
- This gives a measure of the strength of the 2nd degree connection.
- In other words: **how well do you know the mutual friend that can connect you to the 2nd degree node?**



Wedge through Triangle



Recommender Systems: Traditional vs. Graph

Traditional:

- At Scale
- Production Deployments

Graph-Based:

- Exploratory
- Highly Contextual
- Known Rules



Intro

Background

Examples

Our Work

Graph Databases



Applications of Graph Theory

- Social network analysis
- Map / GPS algorithms - shortest distance between two points, etc.
- AI algorithms
- Search engine algorithms



Examples of Graph Applications

- 9/11 Terrorist Network
- London Phone Network
- Enron Emails
- Panama Papers



Intro
Background
Examples

Our Work

Graph Databases



Intro
Background
Examples

Our Work
Graph Databases

Keyword Rec: an HR Keyword Assistant



Consider a typical job requisition

Job Title

- topic information

Company and institute

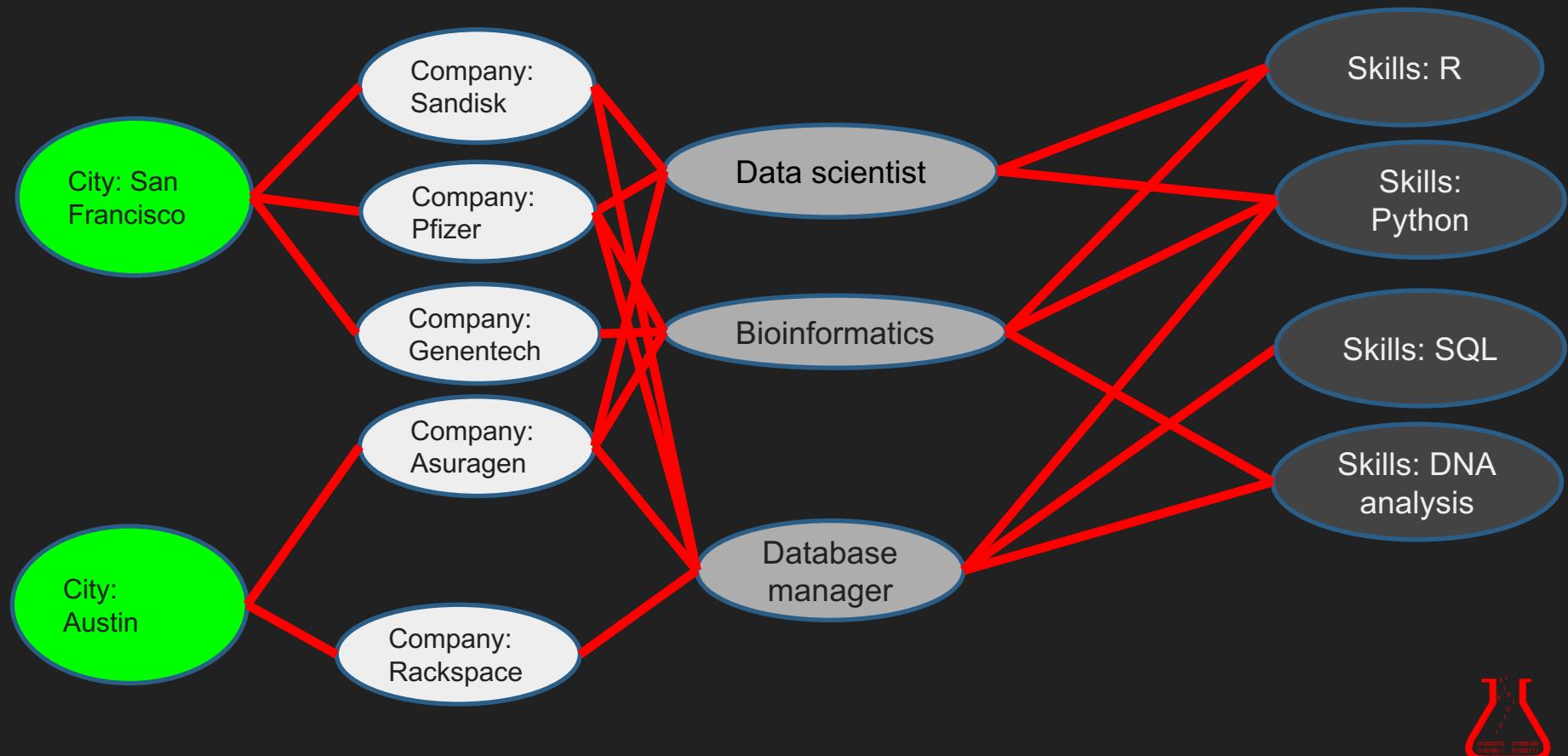
- workplace
- geographic location

Keywords

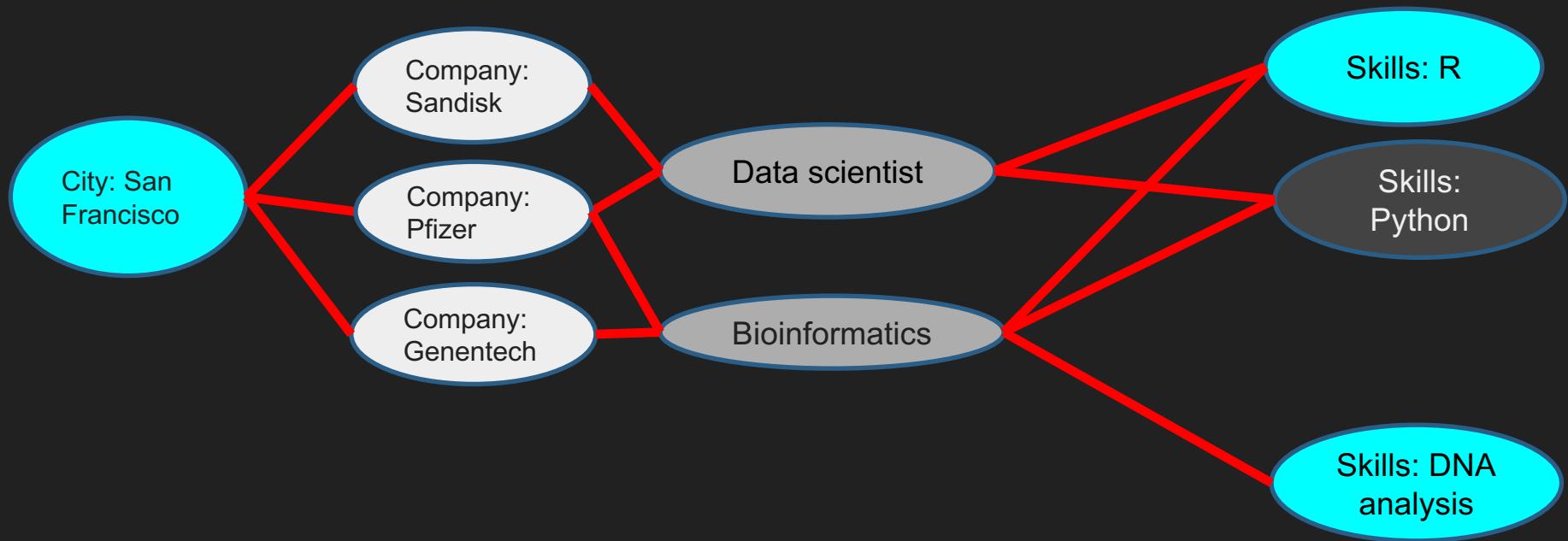
- skills
- experience/focus



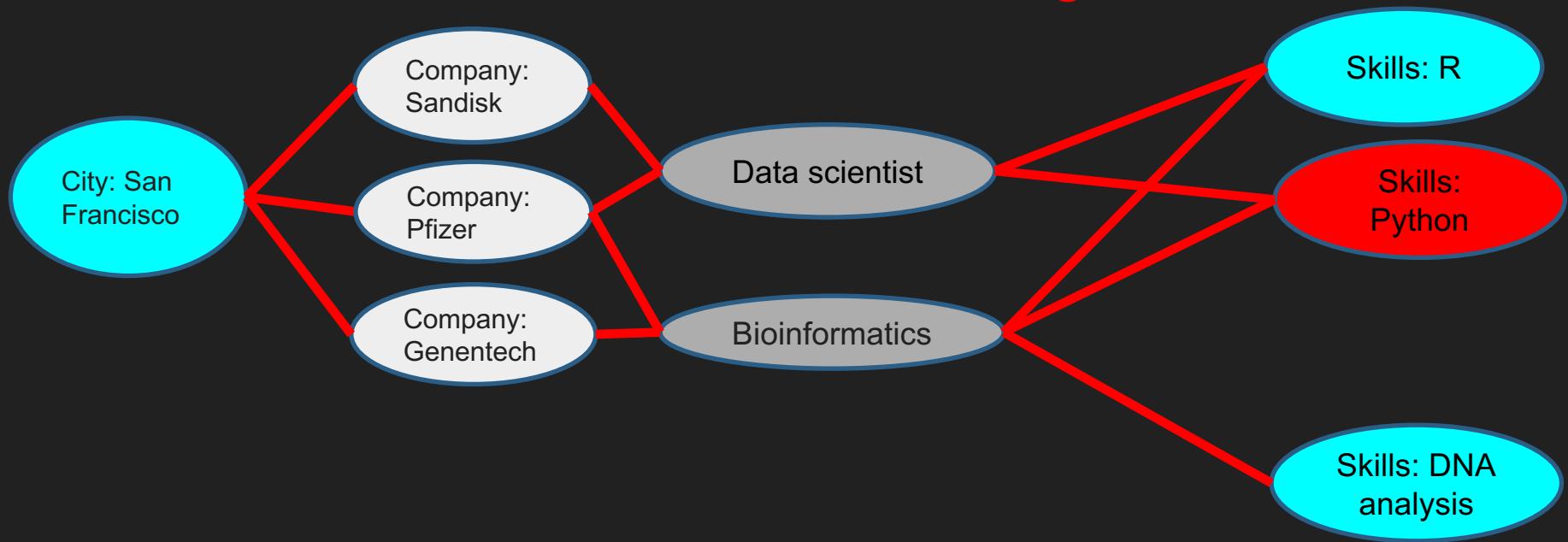
Job requisition analysis by skills



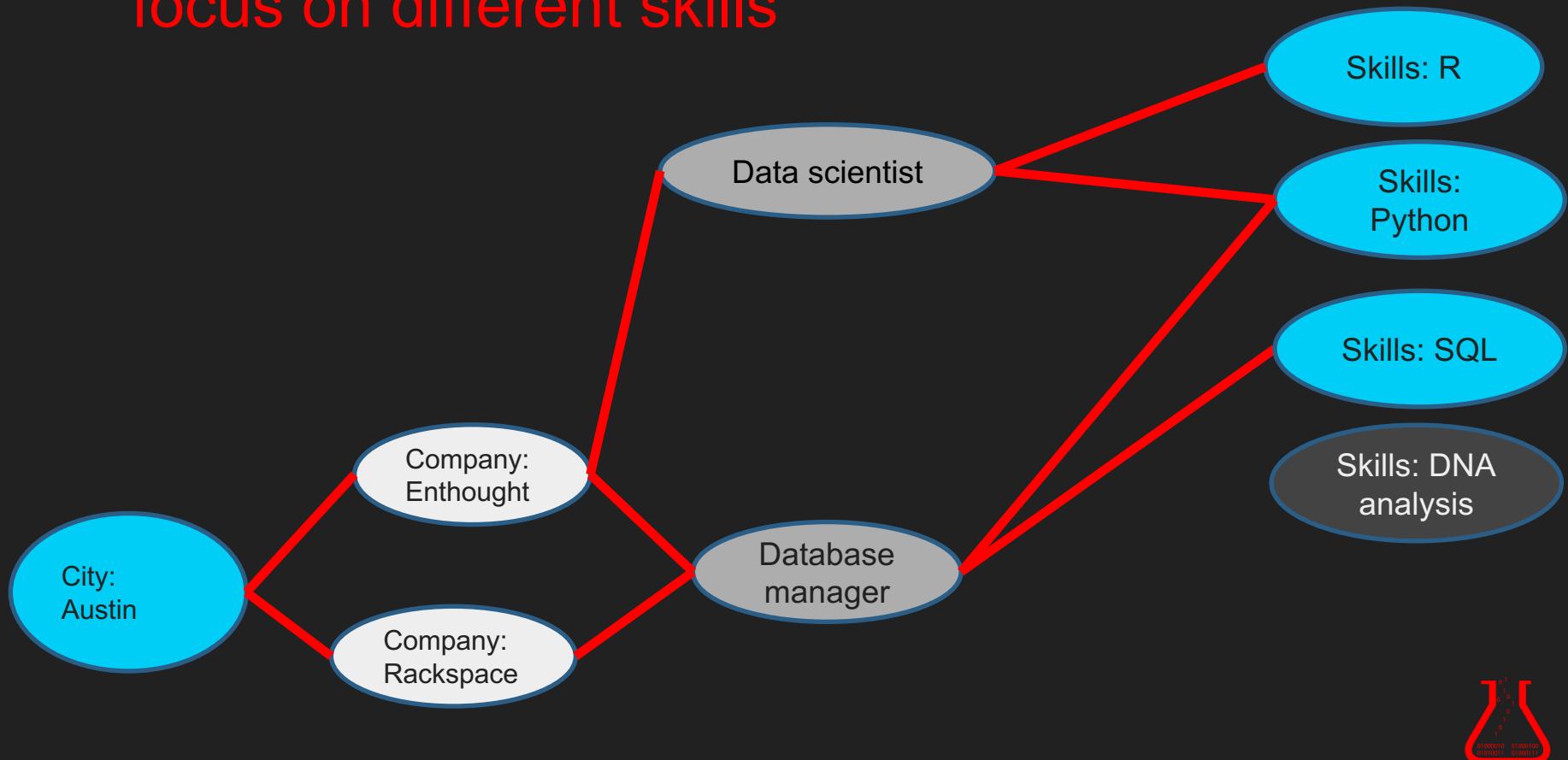
Or looking at a skill might suggest a geographic hub



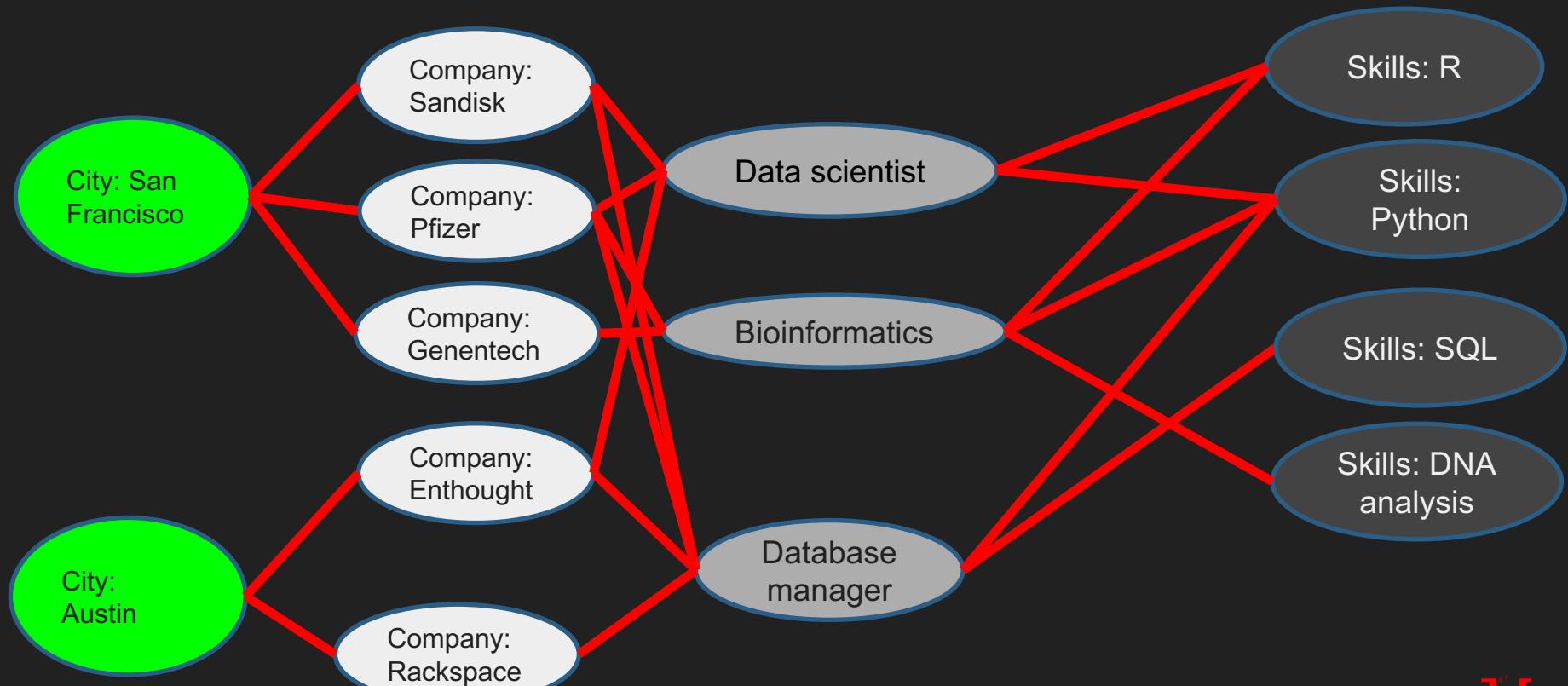
We can also recommend missing skills



If we focus on a specific region we might focus on different skills



Job requisition analysis by skills



The final API



A screenshot of a Jupyter Notebook interface showing three code cells and their outputs.

In [16]: relatedskills('data science','NY',5)

Out[16]: [('data science', 3885), ('big data', 616), ('data engineering', 312), ('science communication', 310), ('information system', 297)]

In [20]: relatedskills('Python','SV',5)

Out[20]: [('python', 328), ('similar', 137), ('unix', 136), ('programming experience', 21), ('language', 15)]

In [21]: relatedskills('Python','NY',5)

Out[21]: [('+', 636), ('programming experience', 615), ('python', 383), ('sql', 324), ('r', 320)]

Intro
Background
Examples

Our Work

Graph Databases

SF vs NY:

Who has better data scientists



Who has Better Data Scientists?

- San Francisco/Silicon Valley
- or New York

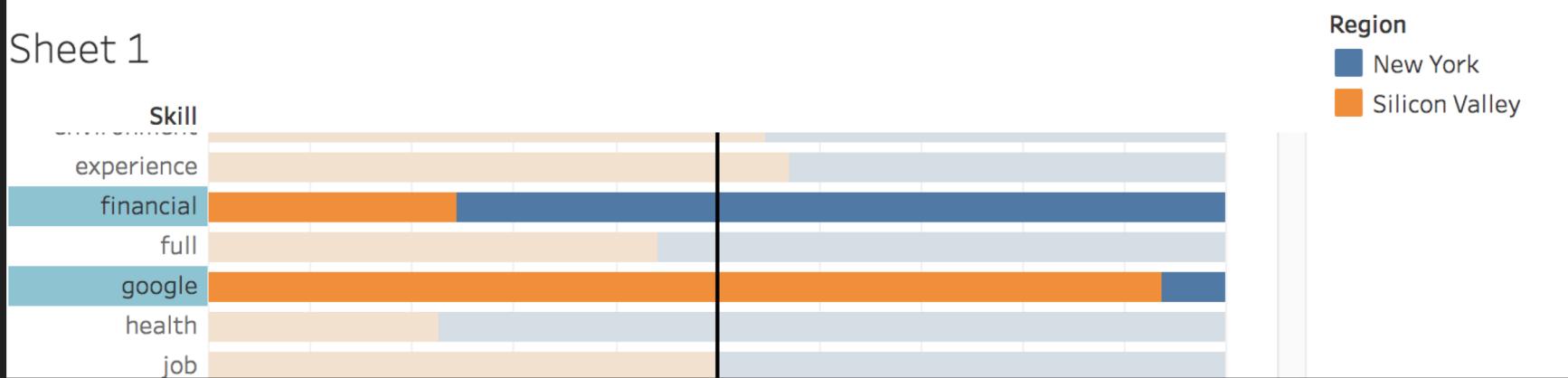
Tableau Link

https://public.tableau.com/profile/denis.vrdoljak#!/vizhome/SVvsNY_JobSkillsPercentCompare/Sheet3?publish=yes

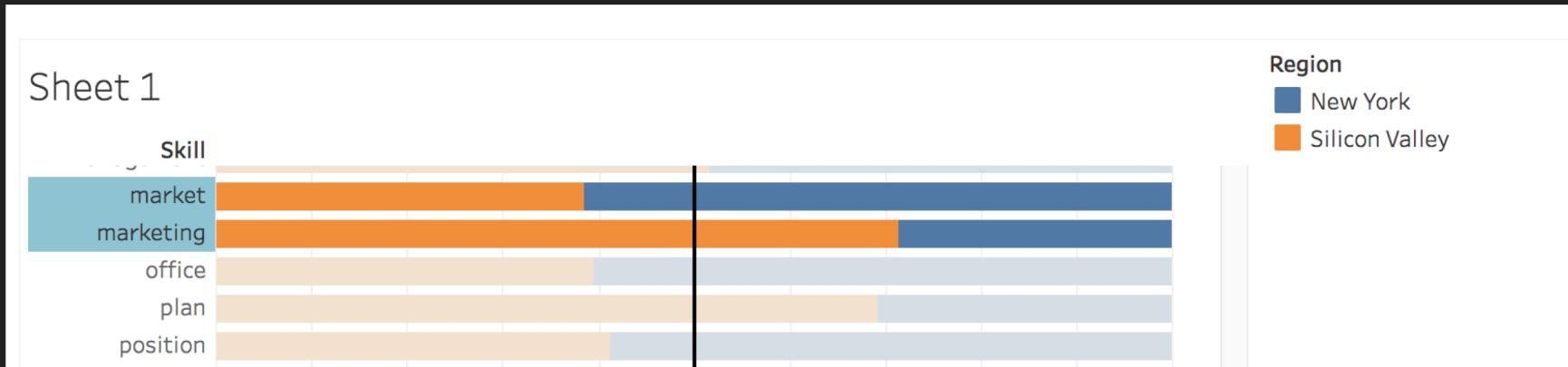


Who has Better Data Scientists?

Sheet 1

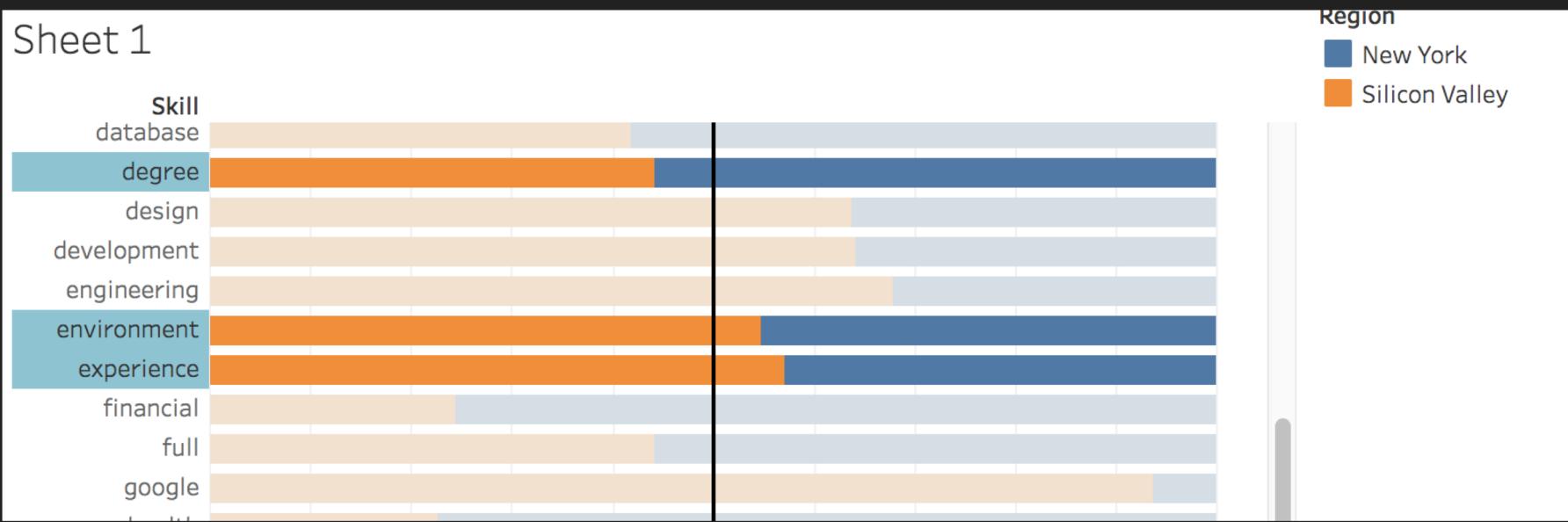


Who has Better Data Scientists?

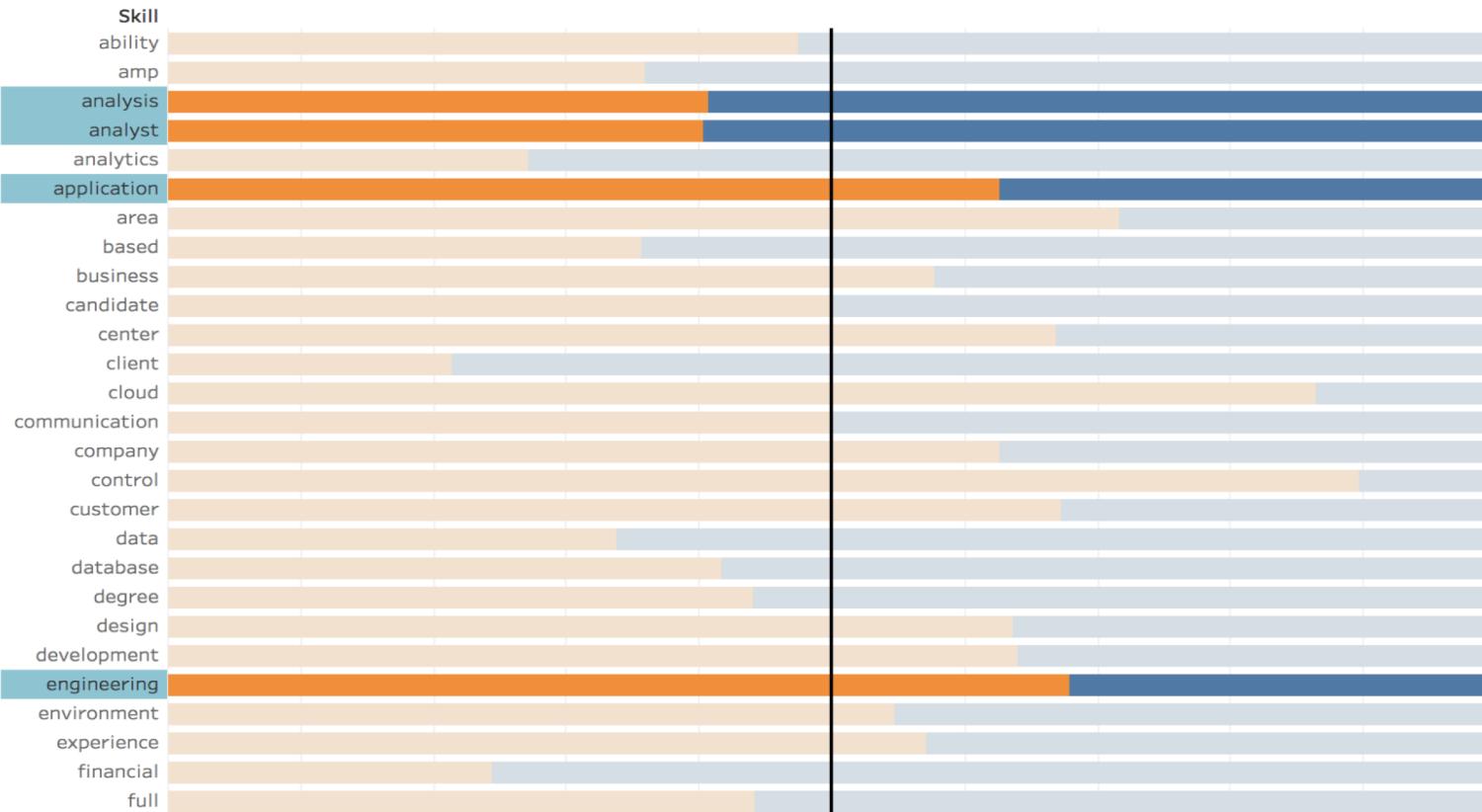


Who has Better Data Scientists?

Sheet 1



Sheet 1



Region

- New York
- Silicon Valley



Intro
Background
Examples

Our Work
Graph Databases

BioRevs:

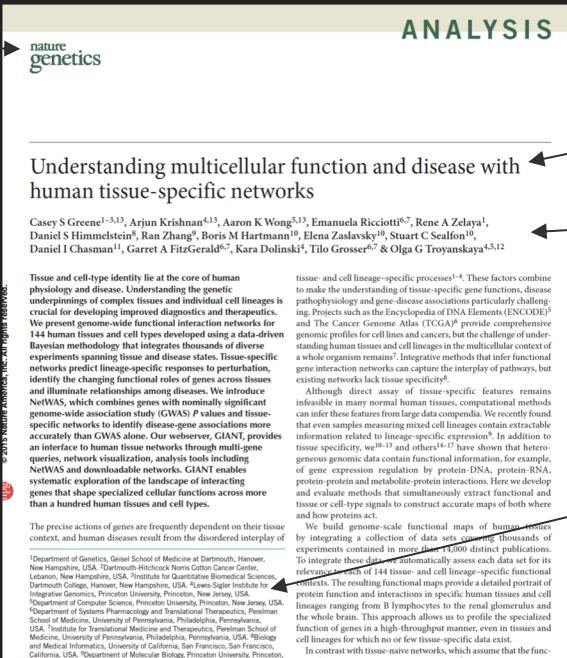
Predicting Biotech IPO Rates
through Collaboration Networks



Publications can be parsed into important data

Journal -

- Discipline
- Audience
- Impact



Title -

- topic information

Collaborator list

- professional relationships

Company and institute

- workplace
- geographic location



Pubmed Database

The world's biomedical research

24,000,000

Full PubMed Dataset (all Biomedical literature)

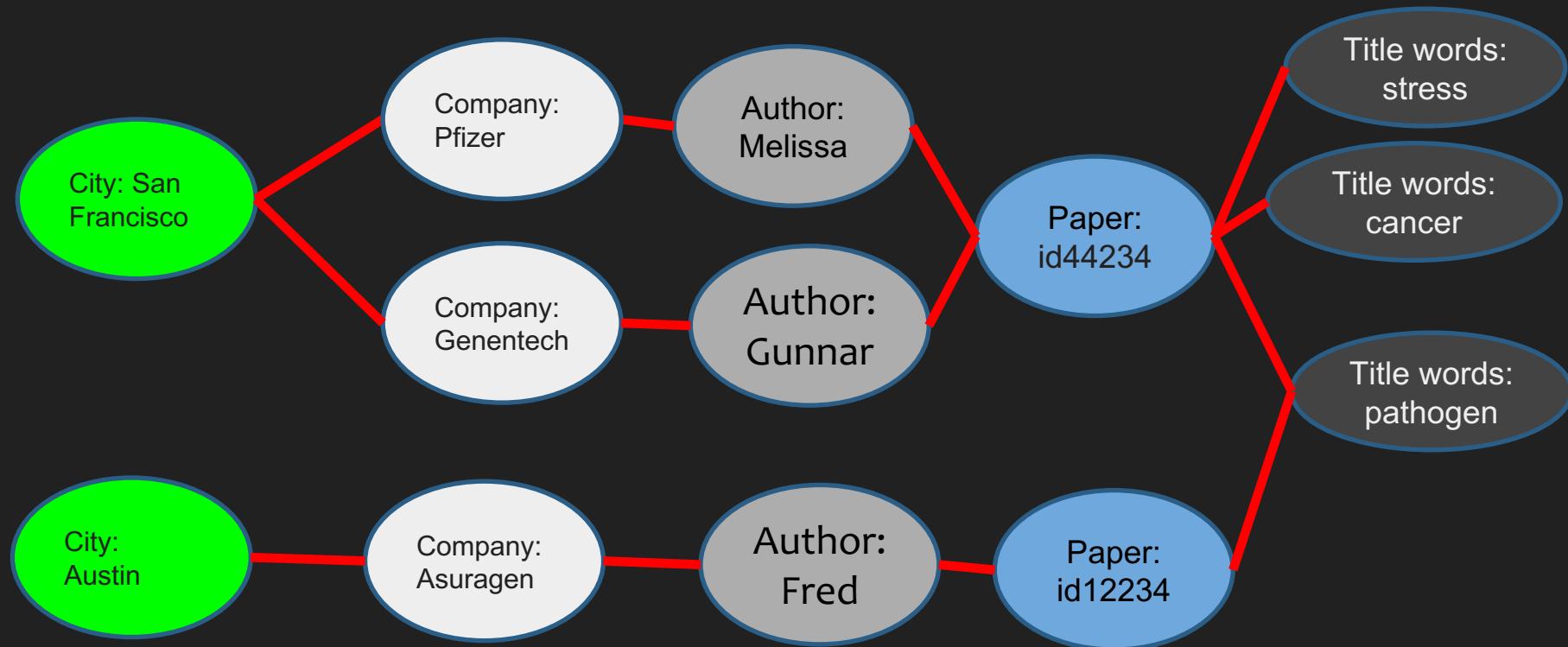
~ 1.5 Million Open Source

PubMed Central (BioTech-open access subset)

Data include: Publication Titles, Authors, and other Metadata



Publication and company data in a graph

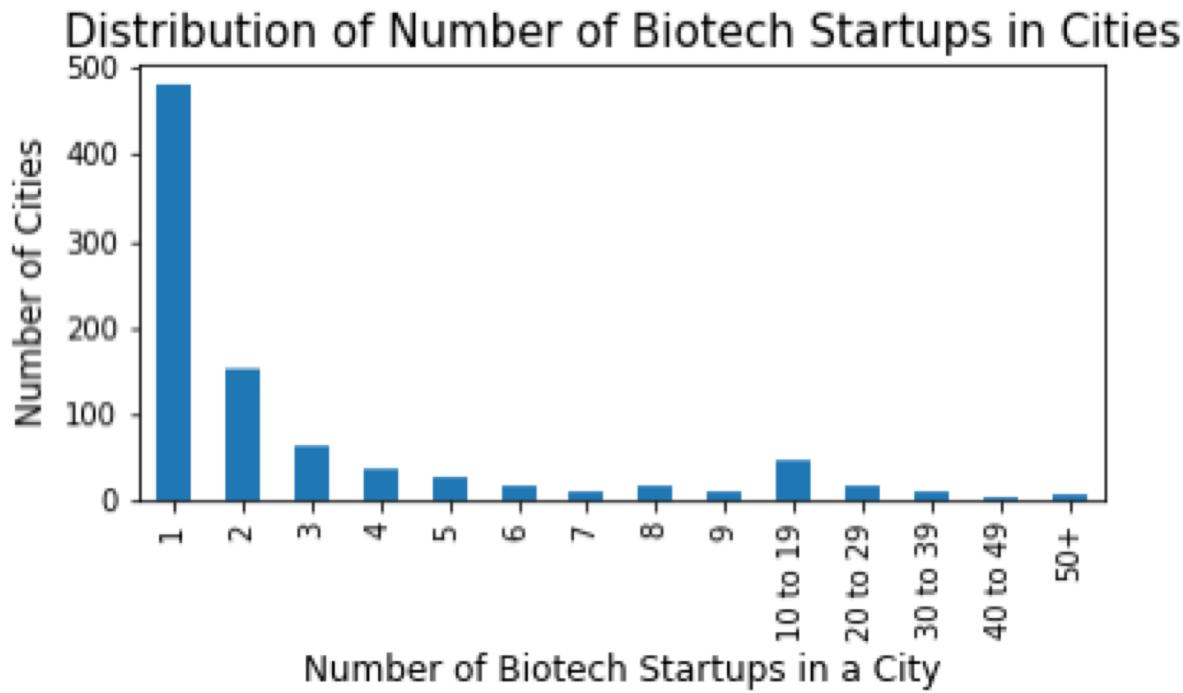


Conducting a job search in biotech

You could identify biotechnology hubs by counting
companies



The Vast Majority of Cities have Few Biotech Startup Companies

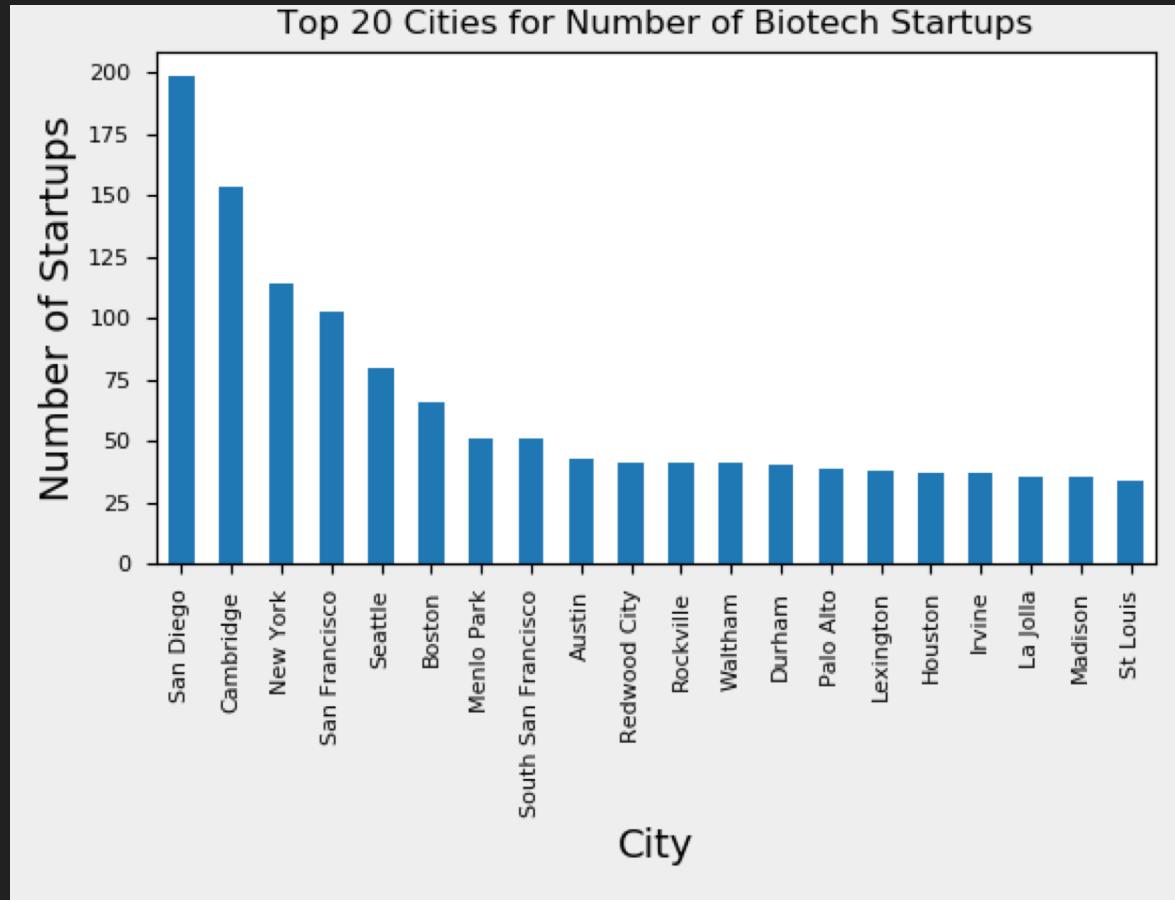


There are 899 US Cities that have at least one Biotech Company.

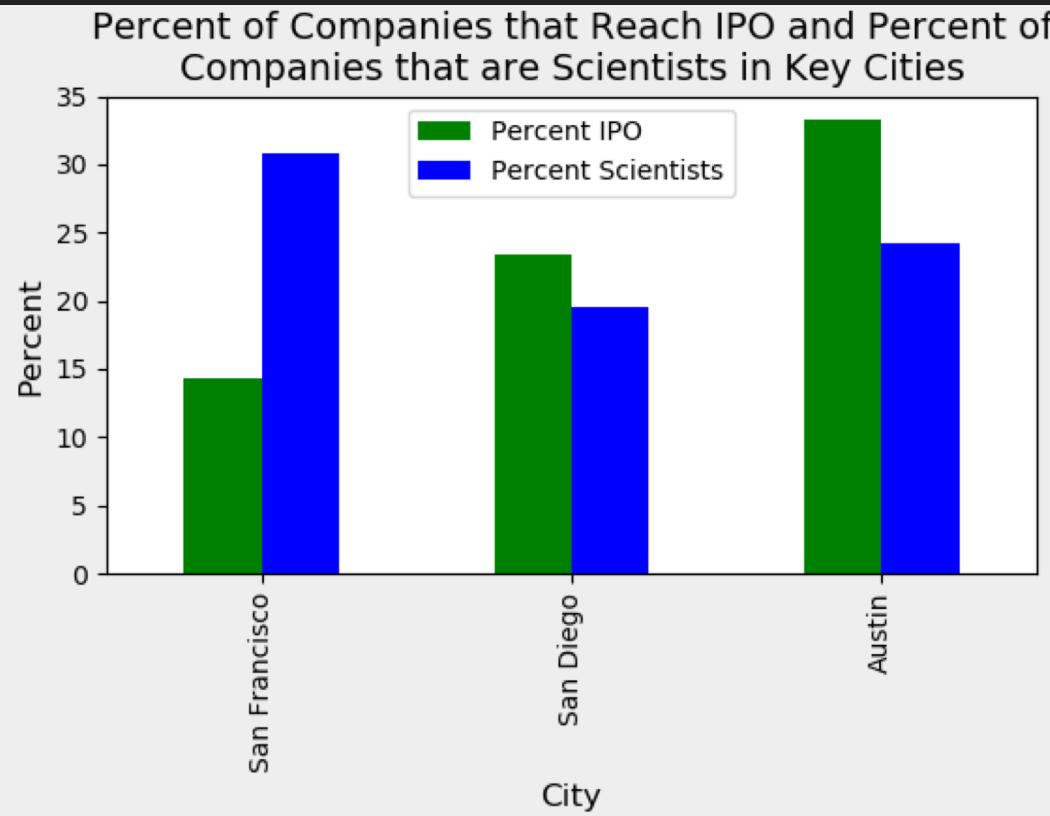
Most cities have very few companies.



We can focus on Some Biotech Hub Cities



Biotech hubs can be profiled by IPO and % scientists



Percent IPO = % Companies reaching IPO in < 6000 days

National average is ~16% reach IPO.

National average is ~21% Scientists.



You could even use public information to examine the
how indexes of company quality

- percent scientists and
- likelihood of success (IPO)

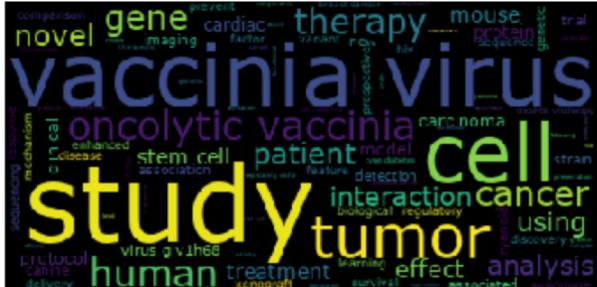


We can use a graph to get more

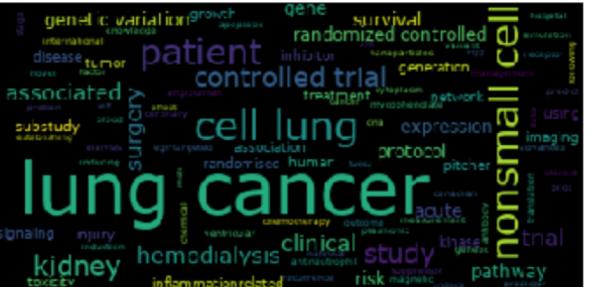


Cities Differ in Scientific Expertise

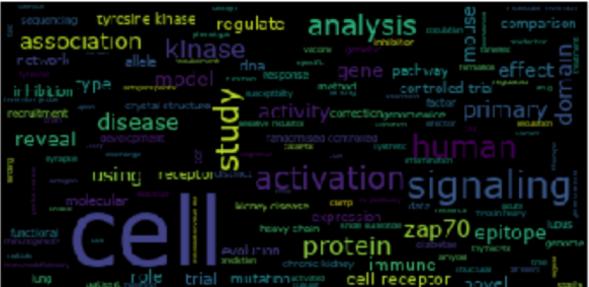
San Diego



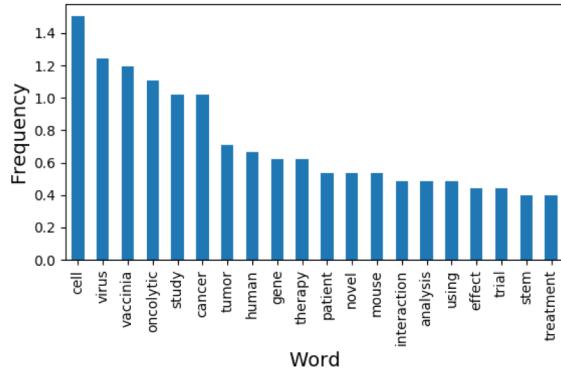
Austin



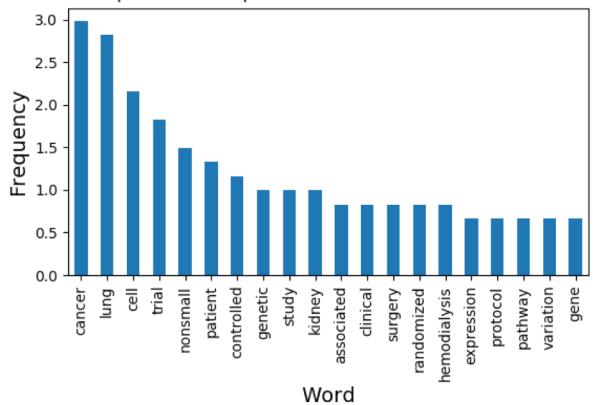
San Francisco



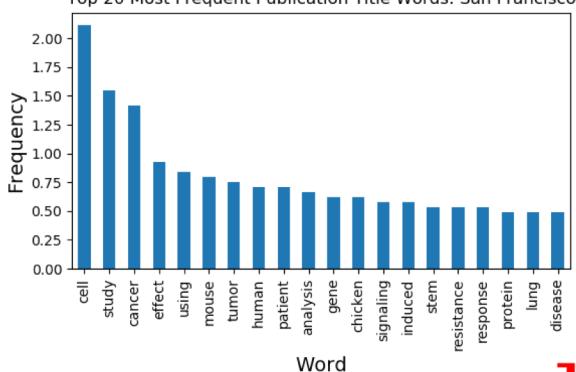
Top 20 Most Frequent Publication Title Words: San Diego



Top 20 Most Frequent Publication Title Words: Austin



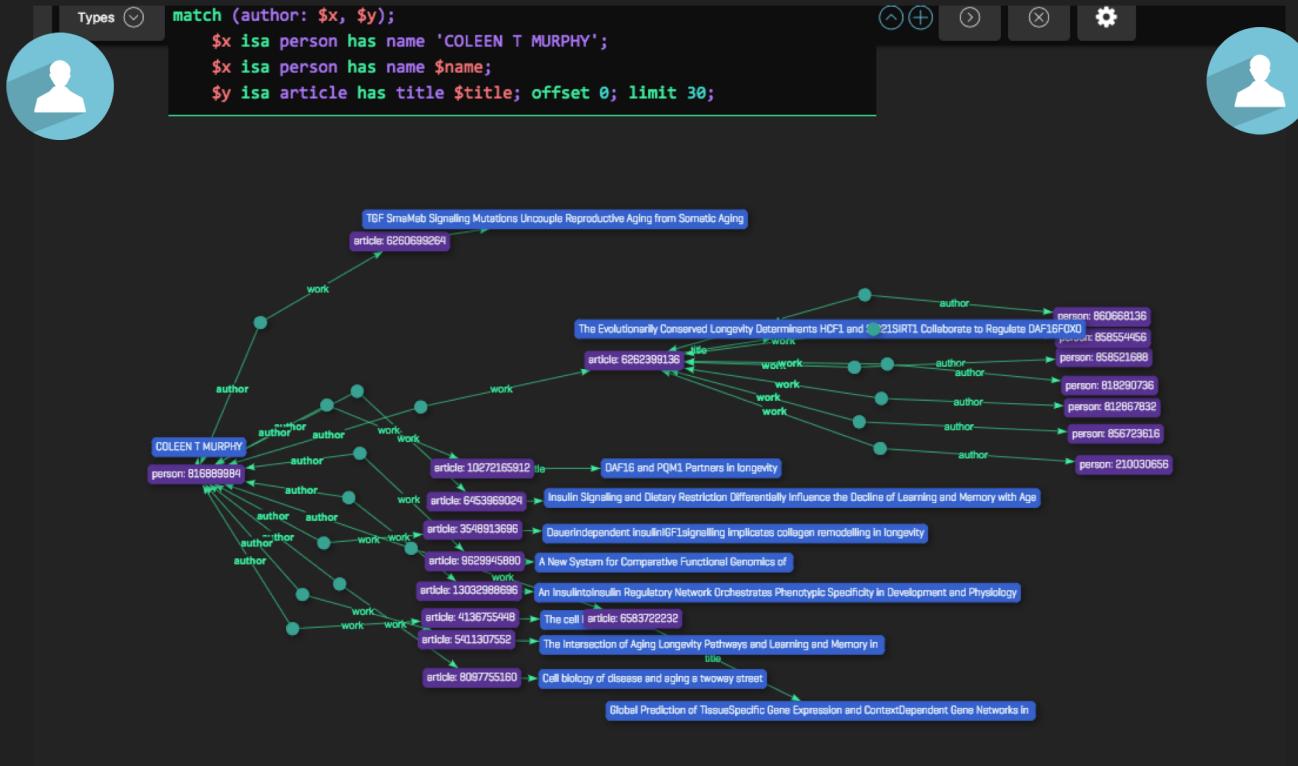
Top 20 Most Frequent Publication Title Words: San Francisco



Quantification of science networks with collaboration graphs



We can get the collaboration network topology from the graph



Finding the top publishers

We identify the top publishers:

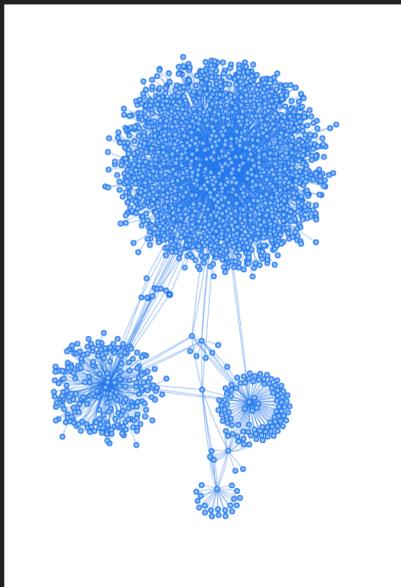
- * Gunther Eysenbach
- * Prof. Hoongkun Fun
- * Seik Weng Ng

```
[neo4j-sh (?)$ MATCH (n)-[r:AUTHORED]->(m)
[>] RETURN n, COUNT(r)
[>] ORDER BY COUNT(r) DESC
[>] LIMIT 10;
+-----+
| n | COUNT(r) |
+-----+
| Node[487153]{name:"THE    "} | 2671 |
| Node[865318]{name:"GUNTHER EYSENBACH"} | 2055 |
| Node[78005]{name:"HOONGKUN FUN"} | 1329 |
| Node[1520956]{name:"SEIK WENG NG"} | 1233 |
| Node[12886]{name:"WEI WANG"} | 1203 |
| Node[33084]{name:"WEI ZHANG"} | 1029 |
| Node[148]{name:"YAN LI"} | 829 |
| Node[6653]{name:"WEI LI"} | 778 |
| Node[15555]{name:"JING WANG"} | 737 |
| Node[20702]{name:"LI ZHANG"} | 670 |
+-----+
10 rows
24526 ms
```

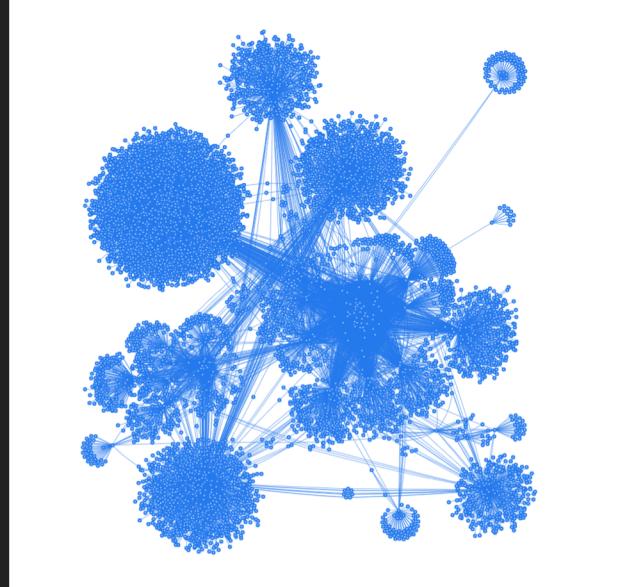


Some typical publication patterns

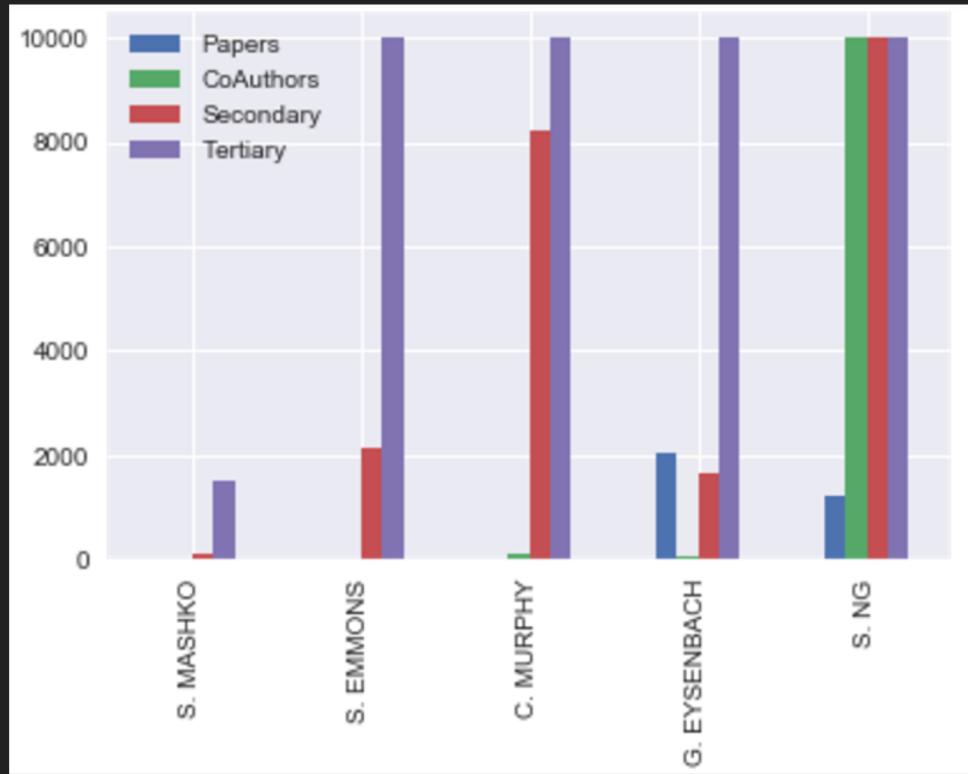
Scott Emmons' 2nd degree collaboration network



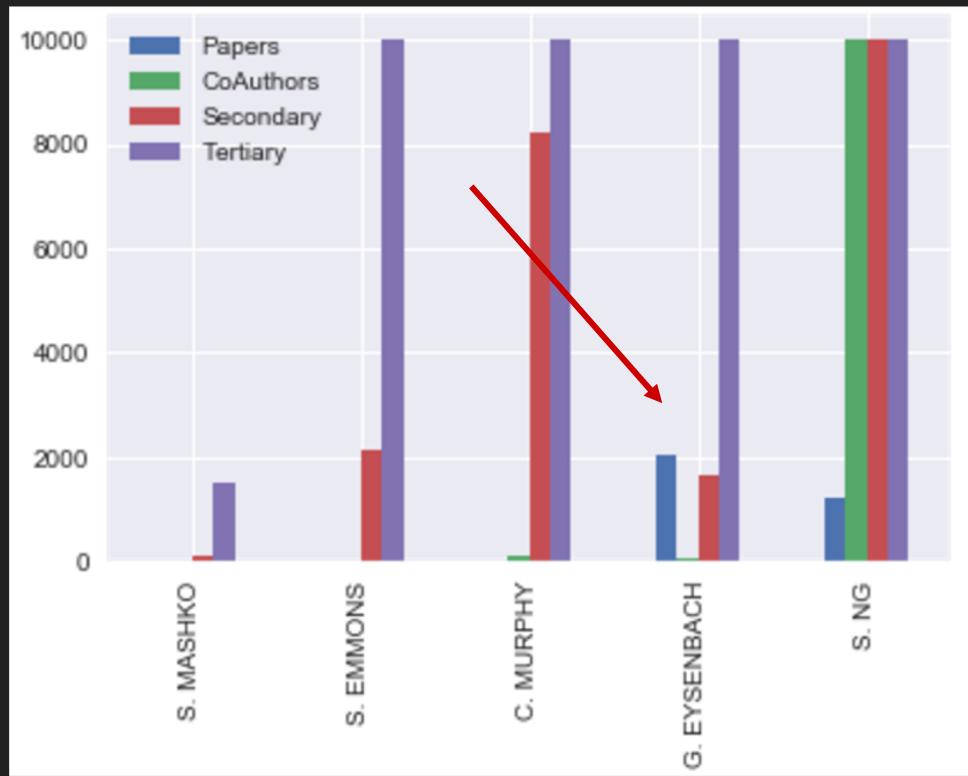
Coleen Murphy's 2nd degree collaboration network



Increasing node count with distance from author



We can see an outlier!



A big publisher with 2055 links!

“one of the most productive researchers, editors, and publishers in the online health field.”

in 2004 received the Janssen-Cilag Future Award, referred to as the German “health care nobel prize”.

Founder of an academic field!
association between search engine queries and influenza incidence,
He coined the terms "infoveillance" and "infodemiology" for these kinds of approaches.



Few co-authors and many papers

Suspicious pattern

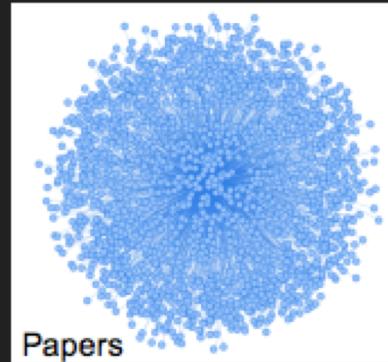
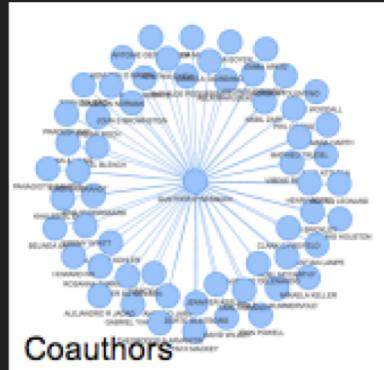
But he is only cited for 120 papers and 40 book chapters.

So what are those other 1900 links ?

Probably editing jobs (false positives)

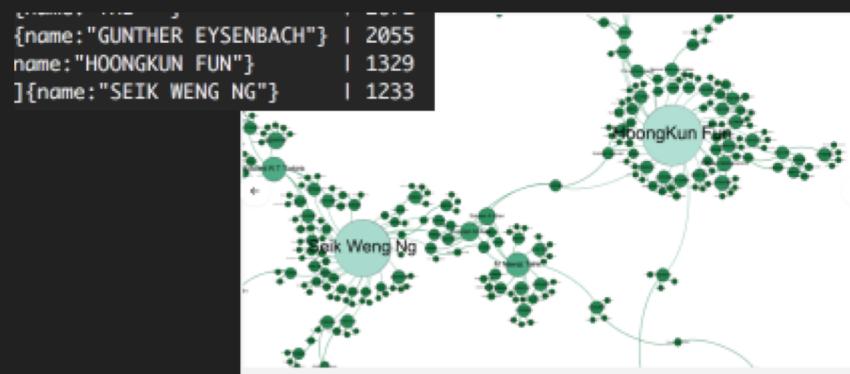
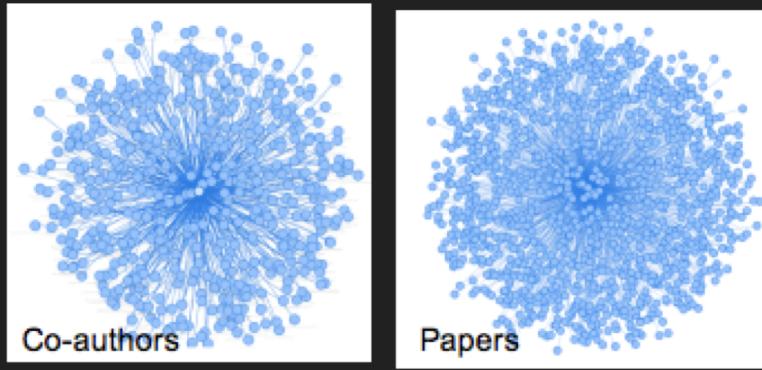


Gunther Eysenbach MUD IMP



Prof. Kun has a pattern closer to what we expected

Many collaborators
and many publications



Intro

Background

Examples

Our Work

Graph Databases



Intro
Background
Examples
Our Work
**Graph
Databases**

Some Key Graph Databases

- Neo4J (well supported)
- Titan/Janus graph (distributed backend)
- AgensGraph (postgres compatible)
- Grakn (knowledge graph)



Neo4J

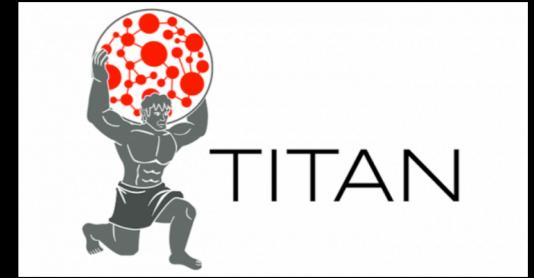


- Industry standard features
- Large userbase and developer community
- Built from the ground up as a graph database

<https://neo4j.com/>



Titan/Janus Graph



- Apache project
- Early adopter of distributed backend
- Elastic scalability
- Integration with Tinkerpop graph stack
- Multiple user access
- Real time updates

<https://www.predictiveanalyticstoday.com/titan/>





AgensGraph

- Highly performant graph database
- Hybrid database built on PostgreSQL
- SQL and Cypher in the same query

<http://bitnine.net/agensgraph/?ckattempt=1>





GRAFN.AI

- Knowledge representation in graphs for AI purposes
 - Nodes represent “objects”, and edges are relationships between them.
- SQL-type query language, Graql, used to quickly and intuitively make queries in the knowledge graph
- Steadily growing technology

<https://grafn.ai/>



Thank You!

Denis Vrdoljak, MIDS
Managing Director, BDSG
Marketing Analyst/Data Scientist, Cisco
denis@bds.group
dvrdolja@cisco.com

Gunnar Kleemann, PhD, MIDS
Senior Data Scientist, BDSG
Data Science Instructor, UC Berkeley
gunnar@bds.group

