

# GRAPH DB'S & APPLICATIONS

DENIS VRDOLJAK | GUNNAR KLEEMANN

UC BERKELEY SCHOOL OF INFORMATION  
BERKELEY DATA SCIENCE GROUP, LLC



# PRES<sup>E</sup>NTATION RO<sup>A</sup>D MAP

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



# Intro

Background

Examples

Our Work

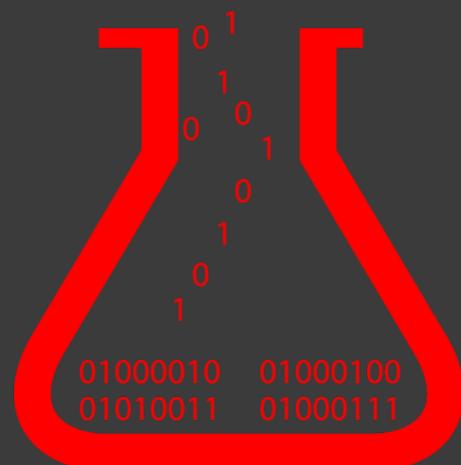
Graph Databases



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](mailto:denis@bds.group)

[dvradolja@cisco.com](mailto:dvradolja@cisco.com)

Gunnar Kleemann, PhD

[gunnar@bds.group](mailto:gunnar@bds.group)

[gunnarkl@berkeley.edu](mailto:gunnarkl@berkeley.edu)



# Why Graphs?

# Why Graph Databases?



# GRAPH DB'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
- GRAPH DB'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

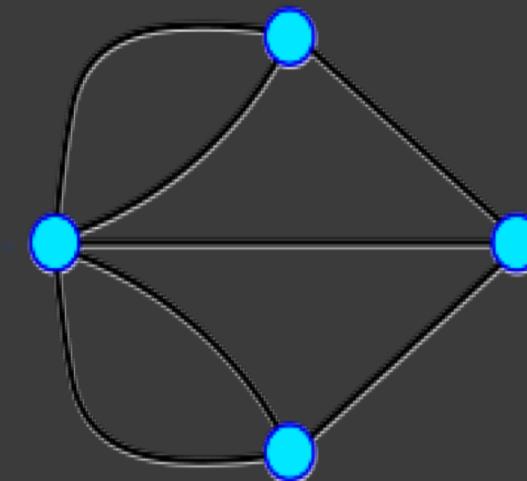
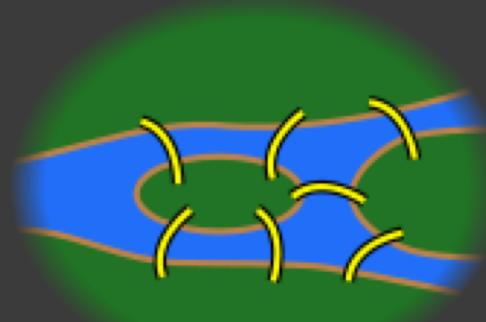
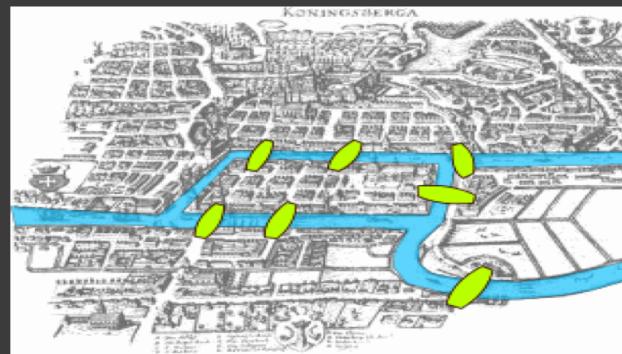
ALSO EVOLVED INTO MODERN DAY NETWORK ANALYSIS (OR NETWORK GRAPH ANALYSIS) AND SOCIAL NETWORK ANALYSIS (SNA)



# THE SEVEN BRIDGES OF KONIGSBERG

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

LEONHARDT EULER CAME UP WITH A NEW WAY OF THINKING ABOUT THE PROBLEM, AND IN TURN BECAME THE FATHER OF MODERN GRAPH THEORY



# 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

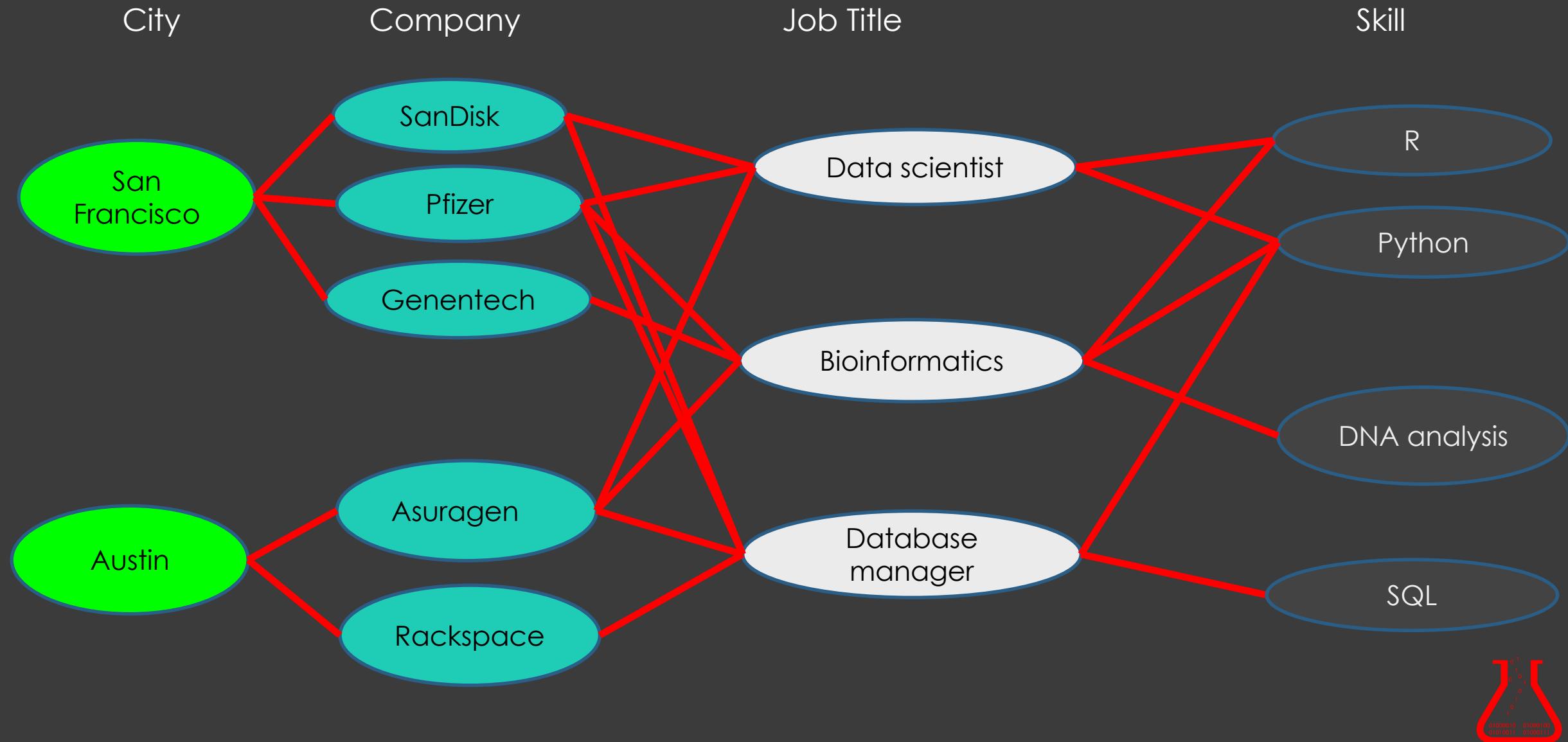
- workplace
- geographic location

## Keywords

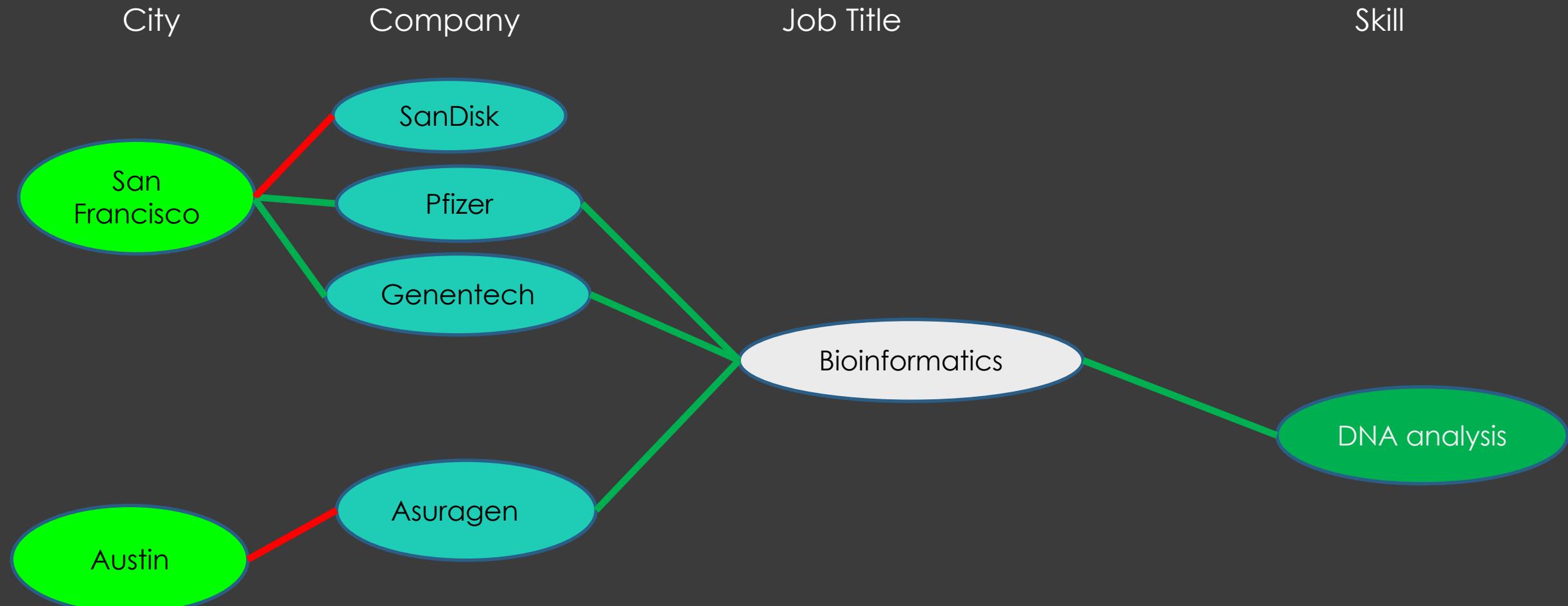
- skills
- experience/focus



# JOB REQUISITION ANALYSIS BY SKILLS



# LOOKING AT A SKILL CAN SUGGEST A GEOGRAPHIC HUB



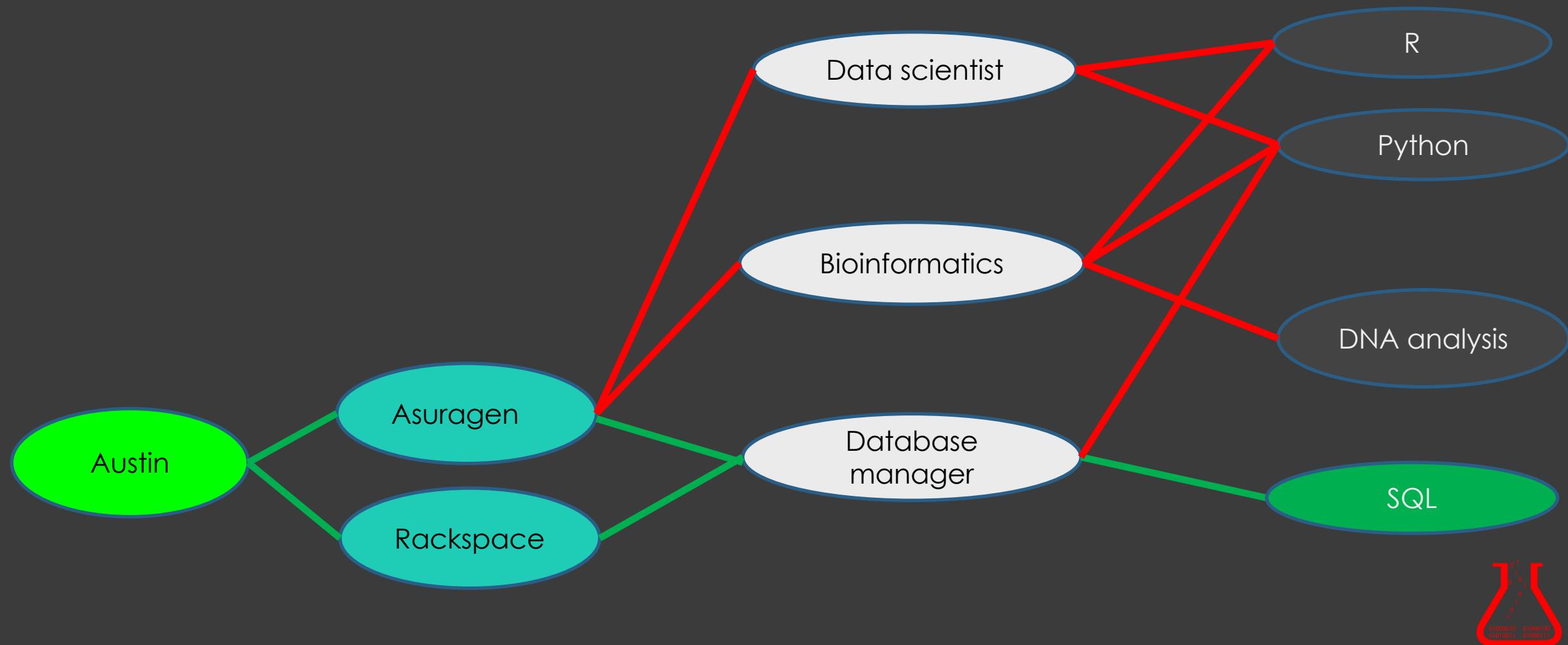
# LOOKING AT A REGION WE MIGHT SUGGEST CRUCIAL SKILLS

City

Company

Job Title

Skill



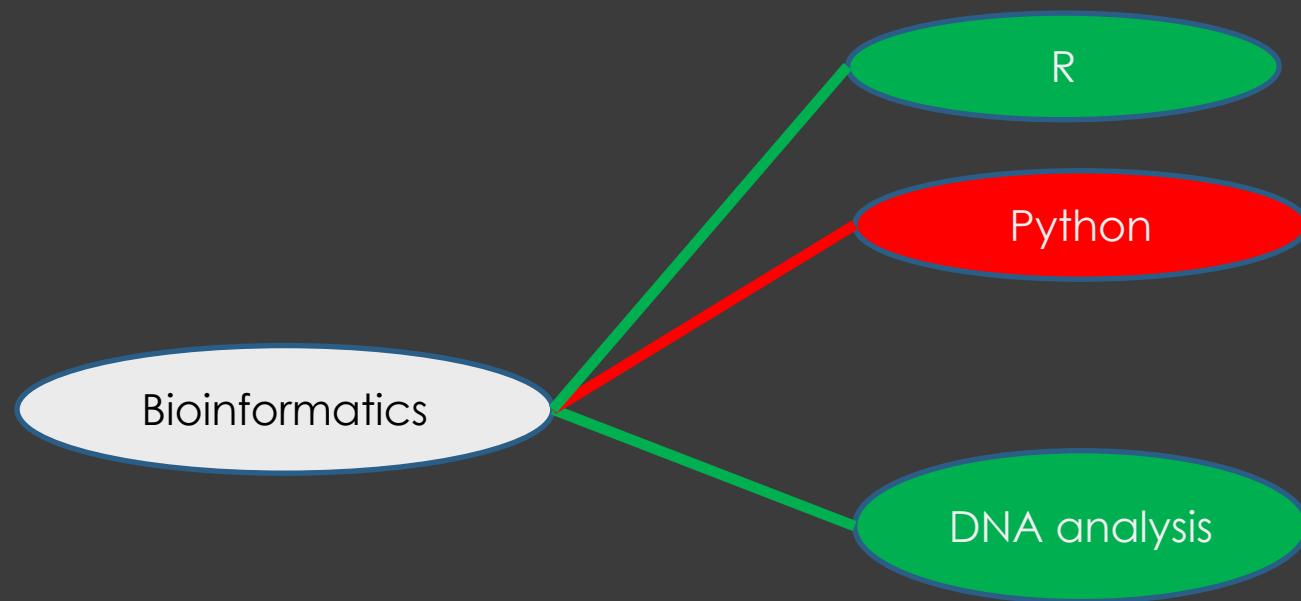
# WE CAN SUGGEST DEFICIENT SKILLS

City

Company

Job Title

Skill



# THE FINAL API

```
relatedskills('Python', 'sv', 5)  
[('python', 328),  
 ('similar', 137),  
 ('unix', 136),  
 ('programming experience', 21),  
 ('language', 15)]
```

```
relatedskills('Python', 'NY', 5)  
[('plus', 636),  
 ('programming experience', 615),  
 ('python', 383),  
 ('sql', 324),  
 ('r', 320)]
```



Intro  
Background  
Examples

## Our Work

Graph Databases

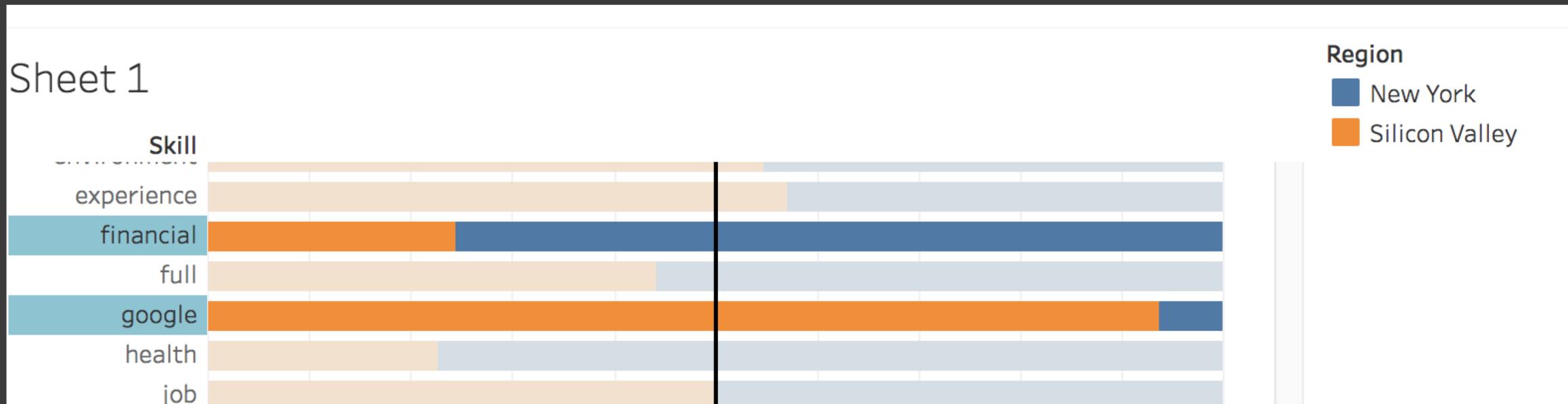
# SF vs NY:

# Who has better data scientists?

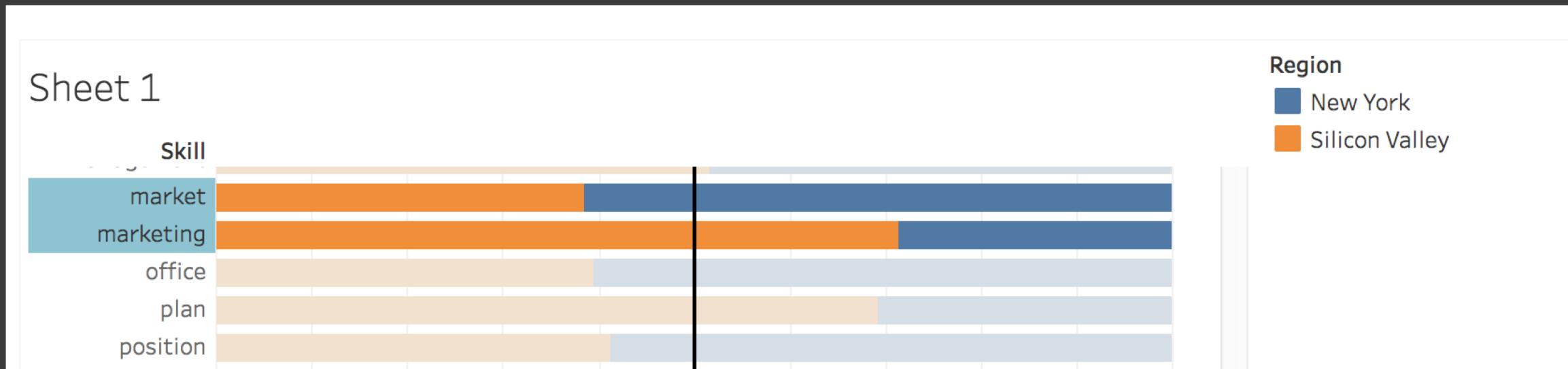
[https://public.tableau.com/profile/denis.vrdoljak#!/vizhome/SVvsNY\\_JobSkillsPercentCompare/Sheet3?publish=yes](https://public.tableau.com/profile/denis.vrdoljak#!/vizhome/SVvsNY_JobSkillsPercentCompare/Sheet3?publish=yes)



# Who has Better Data Scientists?

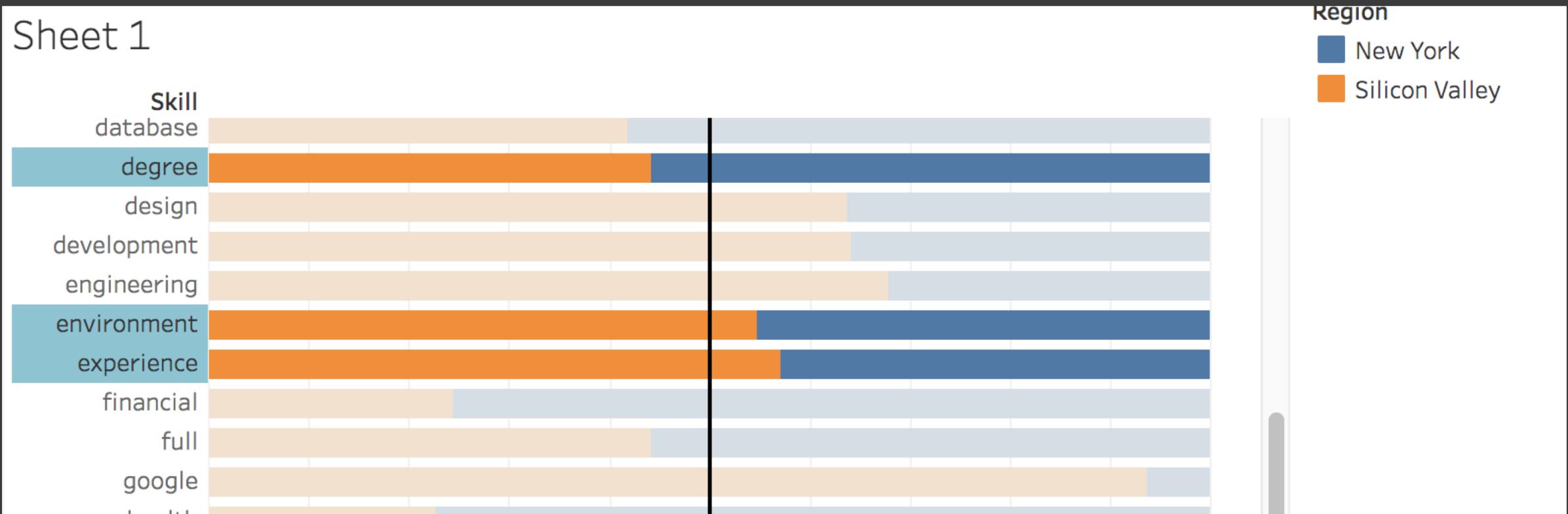


# Who has Better Data Scientists?



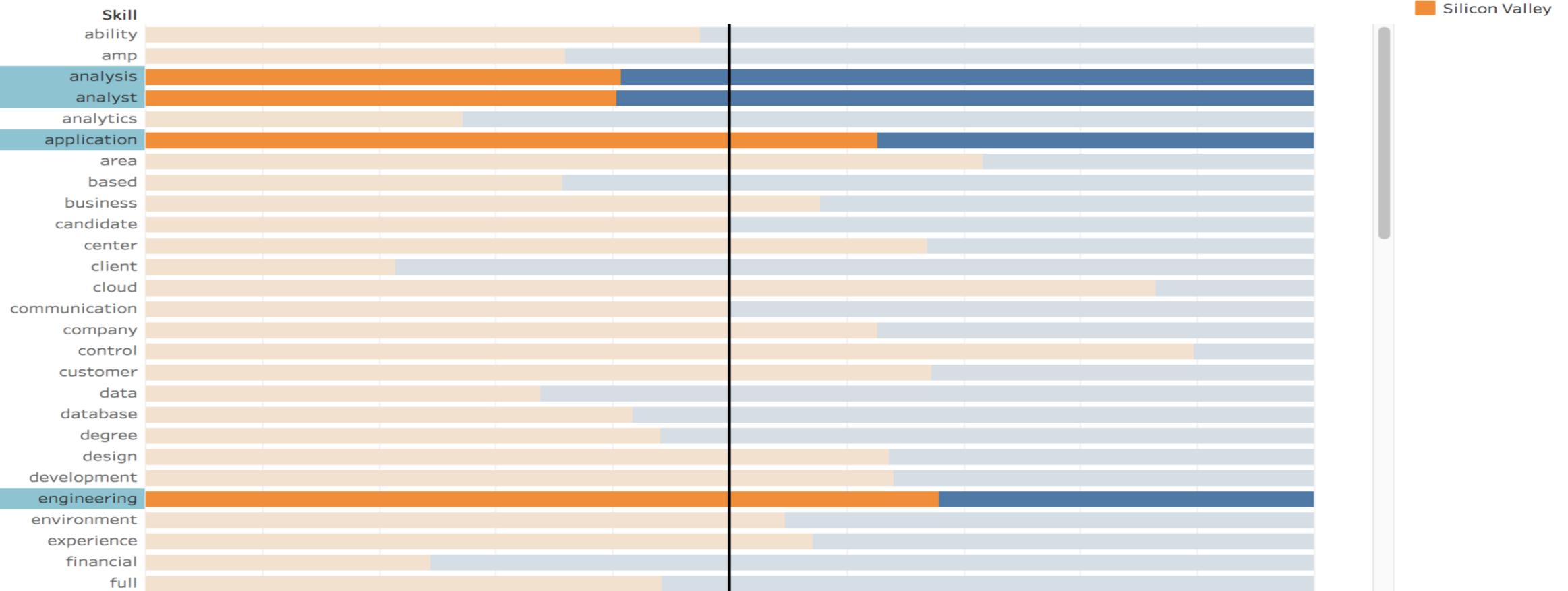
# Who has Better Data Scientists?

Sheet 1



# Who has Better Data Scientists?

Sheet 1





Intro  
Background  
Examples

## Our Work

Graph Databases

# BioRevs:

Predicting Biotech IPO Rates  
through Collaboration Networks

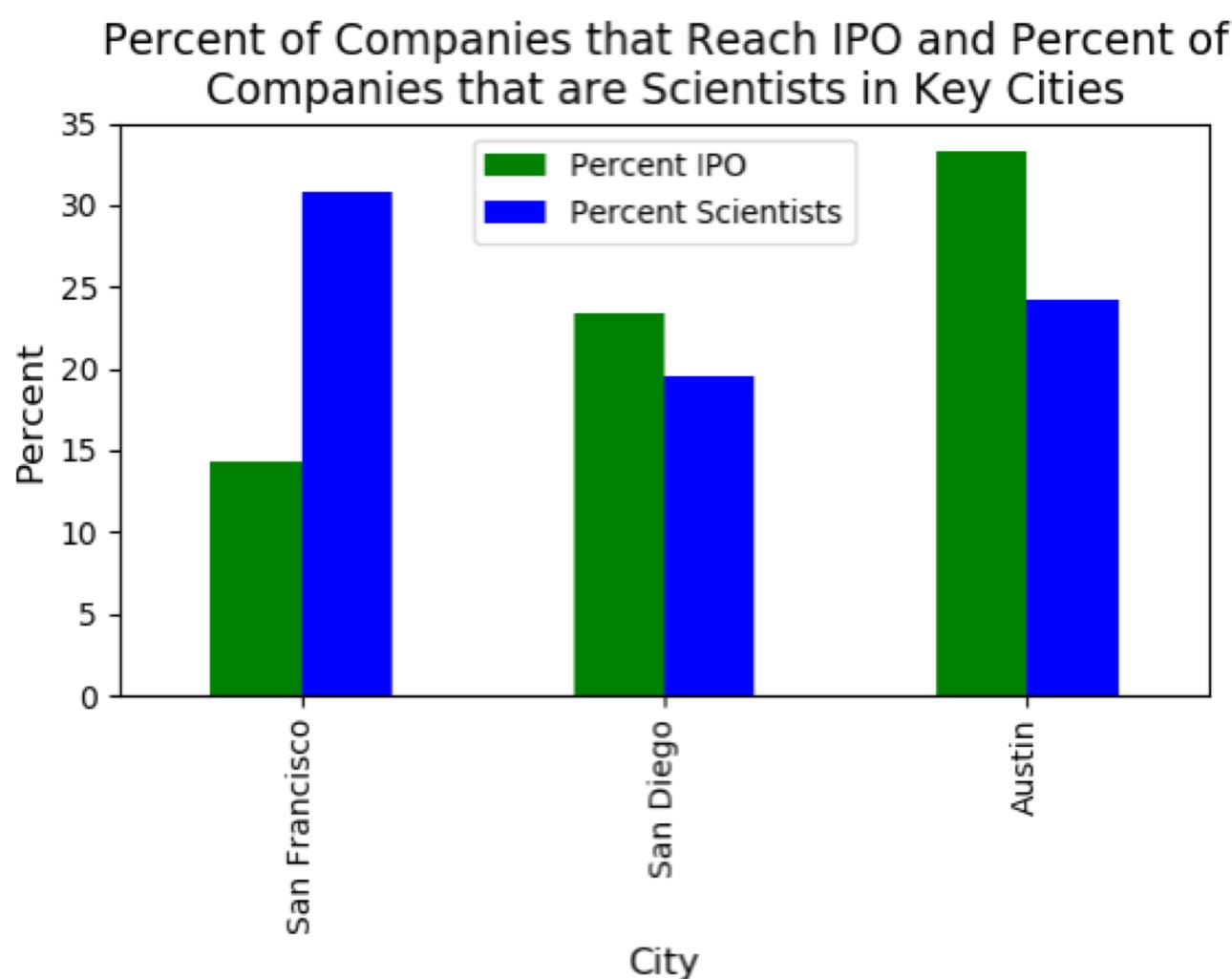


# CONDUCTING A JOB SEARCH IN BIOTECH

YOU COULD IDENTIFY BIOTECHNOLOGY HUBS BY  
COUNTING COMPANIES



# BIOTECH HUBS CAN BE PROFILED BY IPO AND % SCIENTISTS



Percent IPO = % Companies reaching IPO in < 6000 days

National average is ~16% reach IPO.

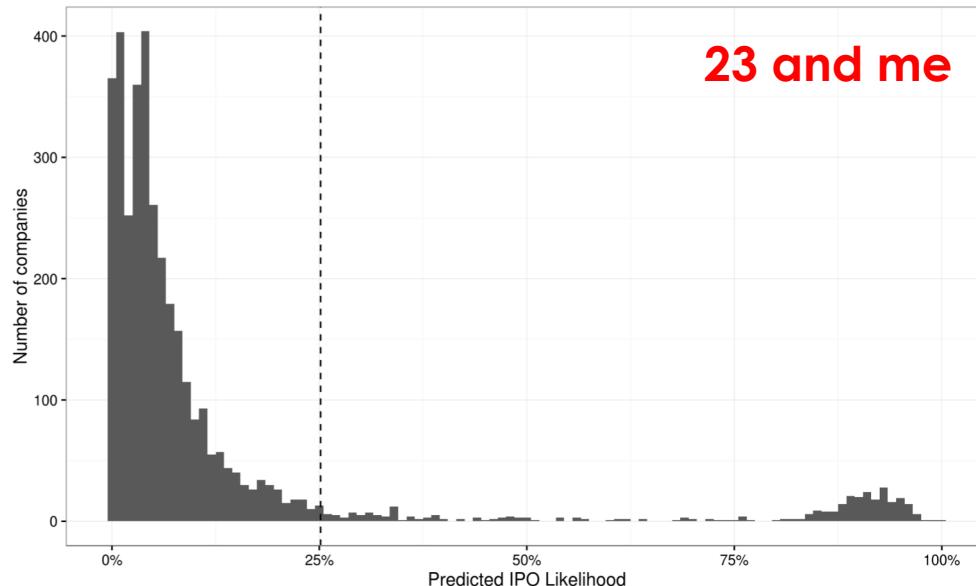
National average is ~21% Scientists.



# BASED ON TRADITIONAL MACHINE LEARNING ANALYSIS WE BUILT BIOREVS

## Success Likelihood

IPO Likelihood: 25.2%

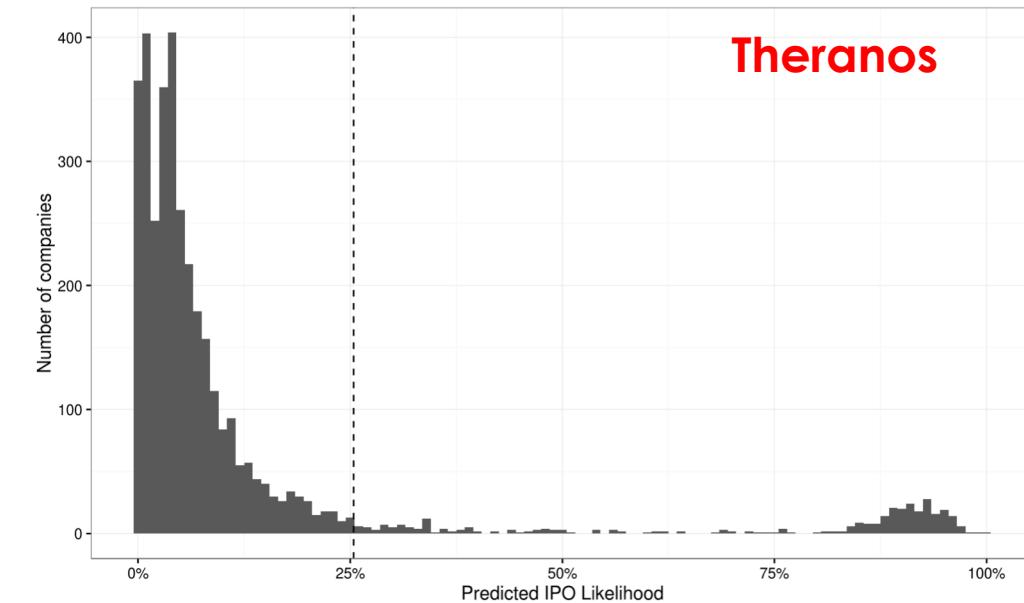


## Expected Valuation

Forecasted IPO Amount: \$521,812,444

## Success Likelihood

IPO Likelihood: 25.4%



## Expected Valuation

Forecasted IPO Amount: \$133,861,306

\* Data up to 2013



WE CAN USE A GRAPH TO  
GET MORE



# 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)

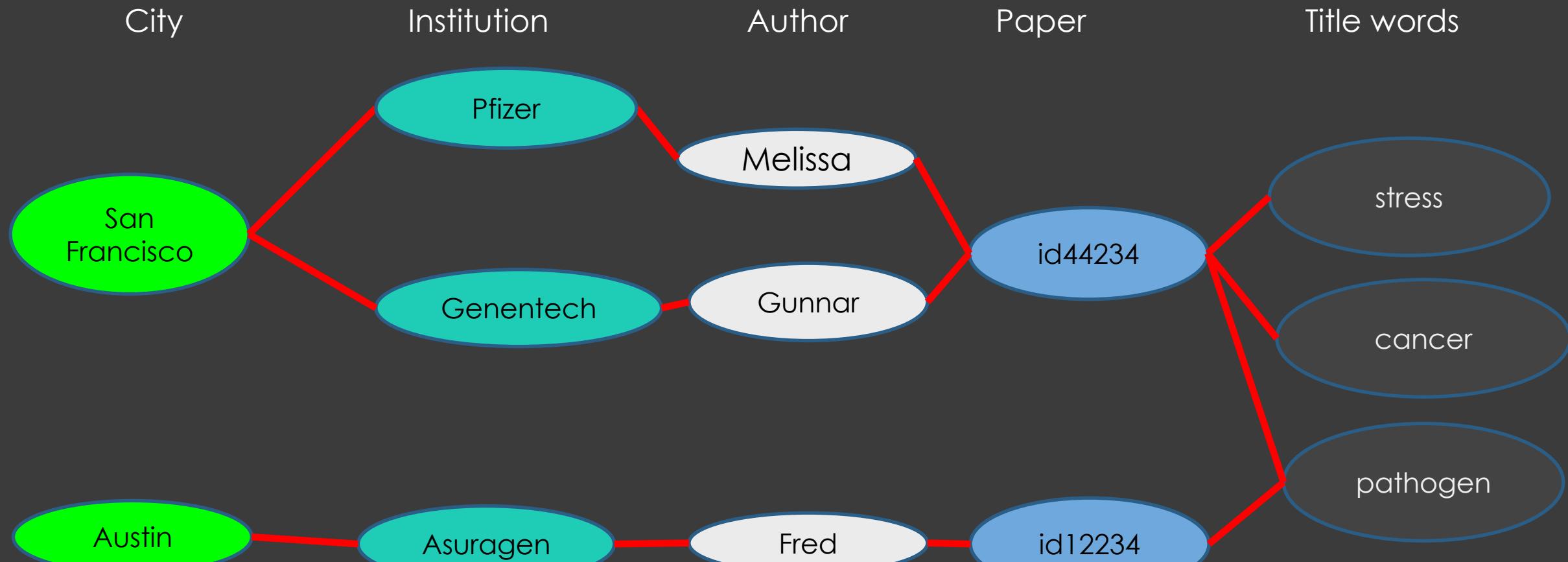
**~ 1.5 MILLION OPEN SOURCE**

PUBMED CENTRAL (BIOTECH-OPEN ACCESS SUBSET)

DATA INCLUDE: PUBLICATION TITLES, AUTHORS, AND OTHER METADATA

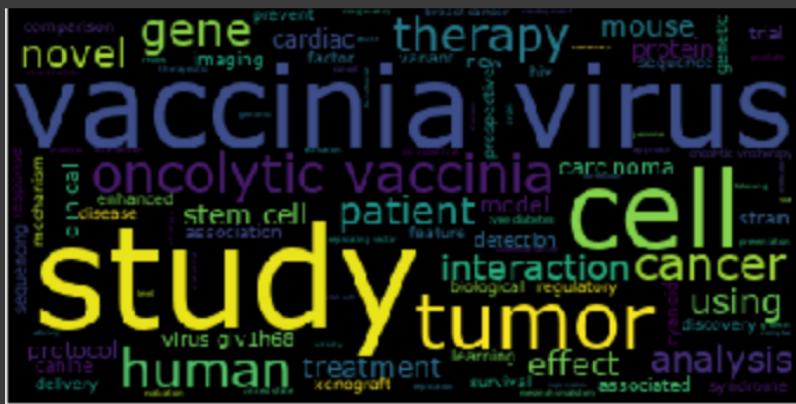


# PUBLICATION AND COMPANY DATA IN A GRAPH

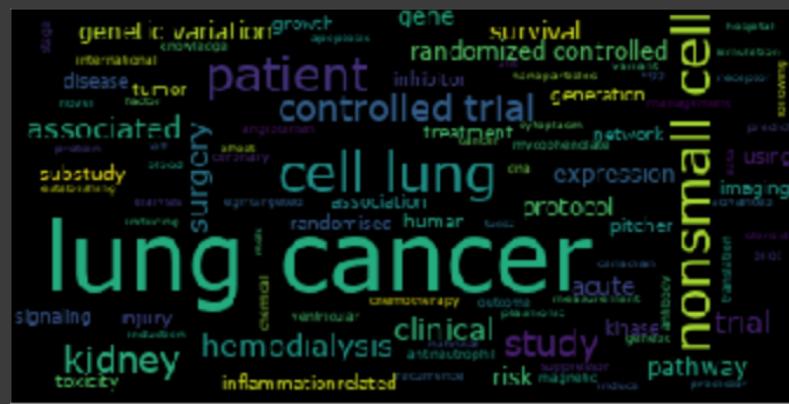


# CITIES DIFFER IN SCIENTIFIC EXPERTISE

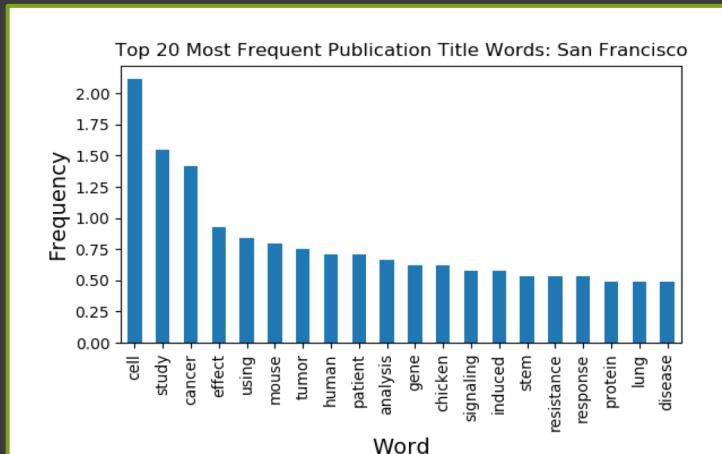
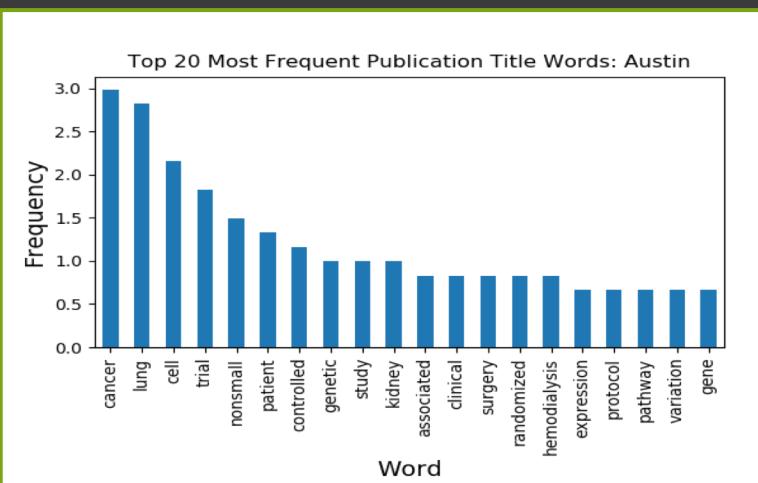
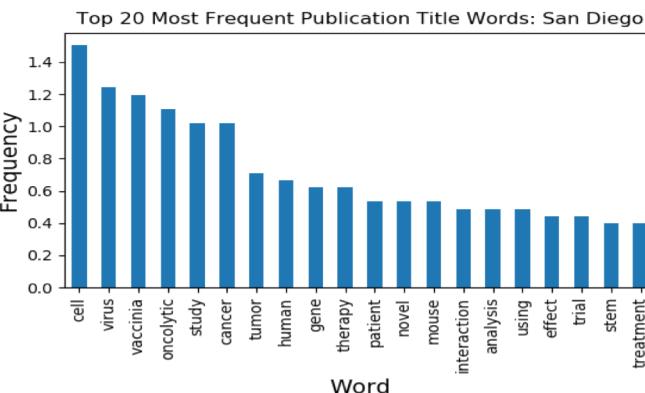
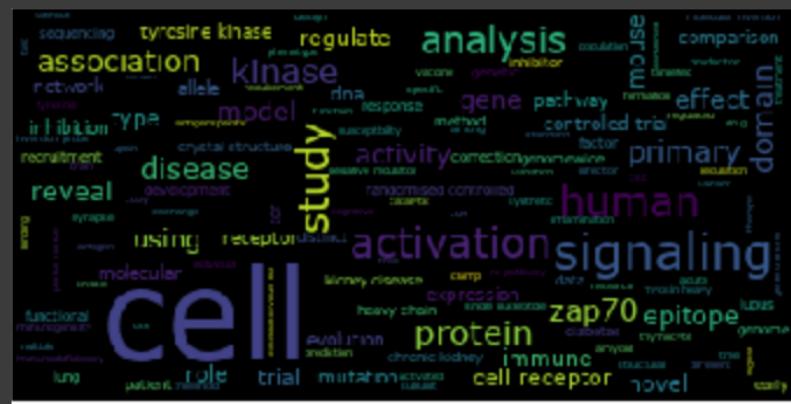
San Diego



Austin



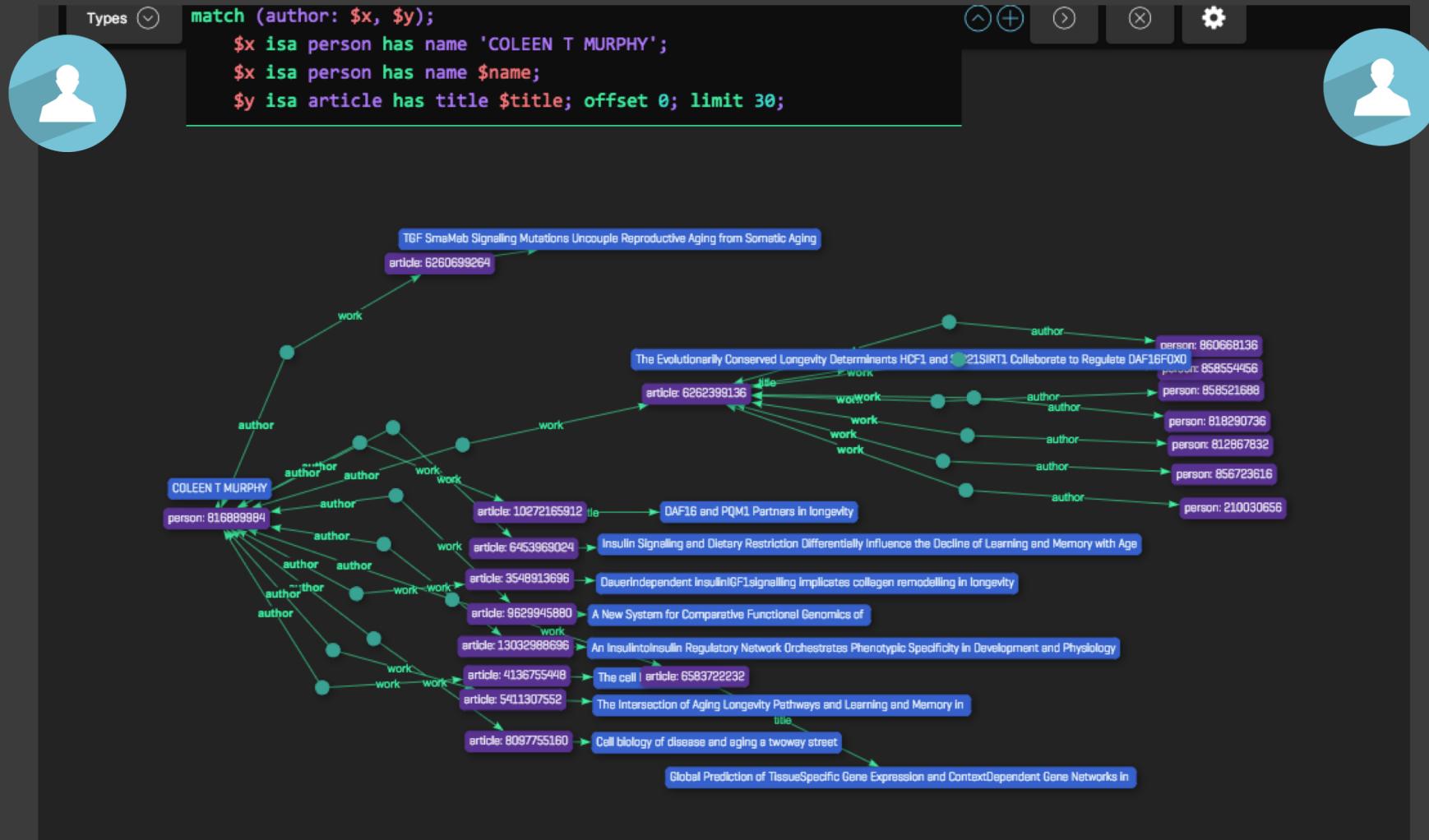
San Francisco



# QUANTIFICATION OF SCIENCE NETWORKS WITH COLLABORATION GRAPHS

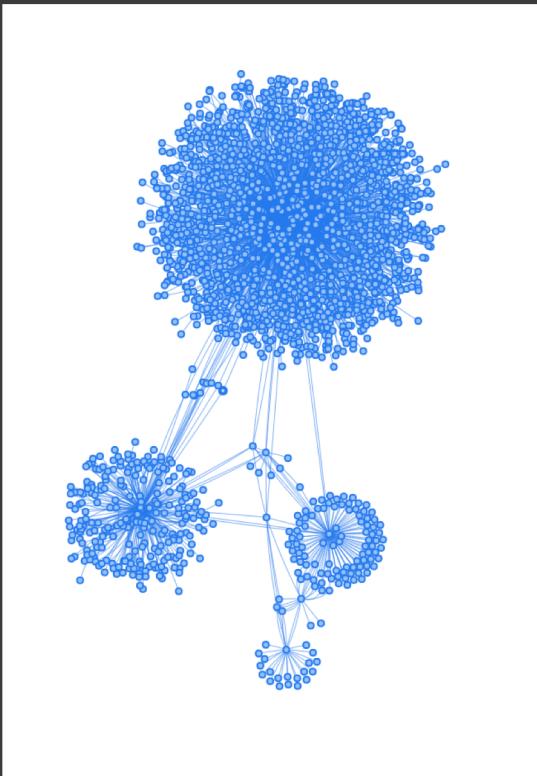


# WE CAN GET THE COLLABORATION NETWORK TOPOLOGY FROM THE GRAPH

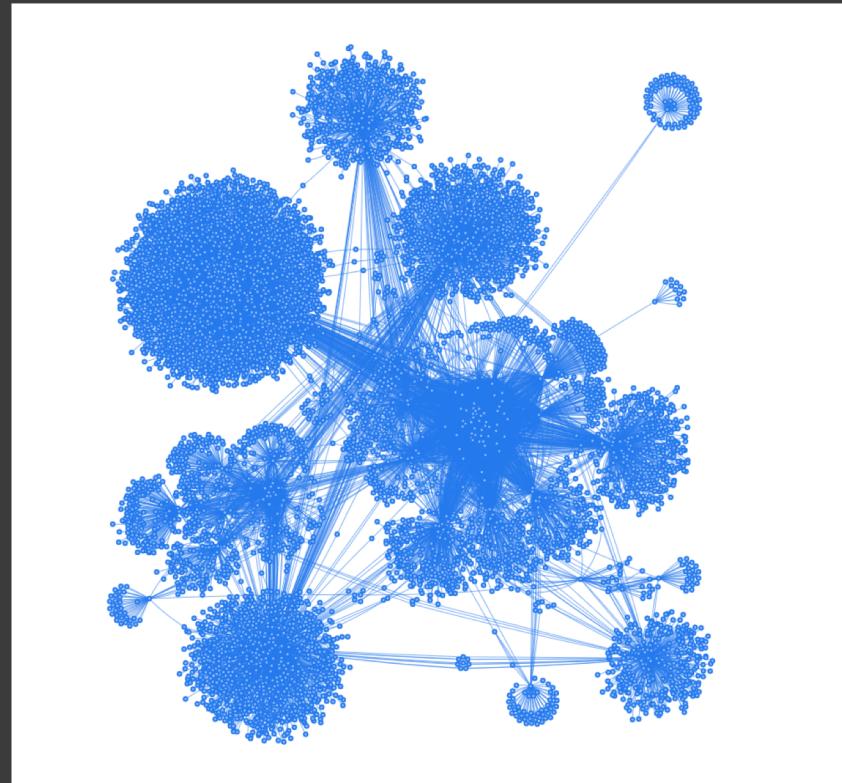


# SOME TYPICAL PUBLICATION PATTERNS

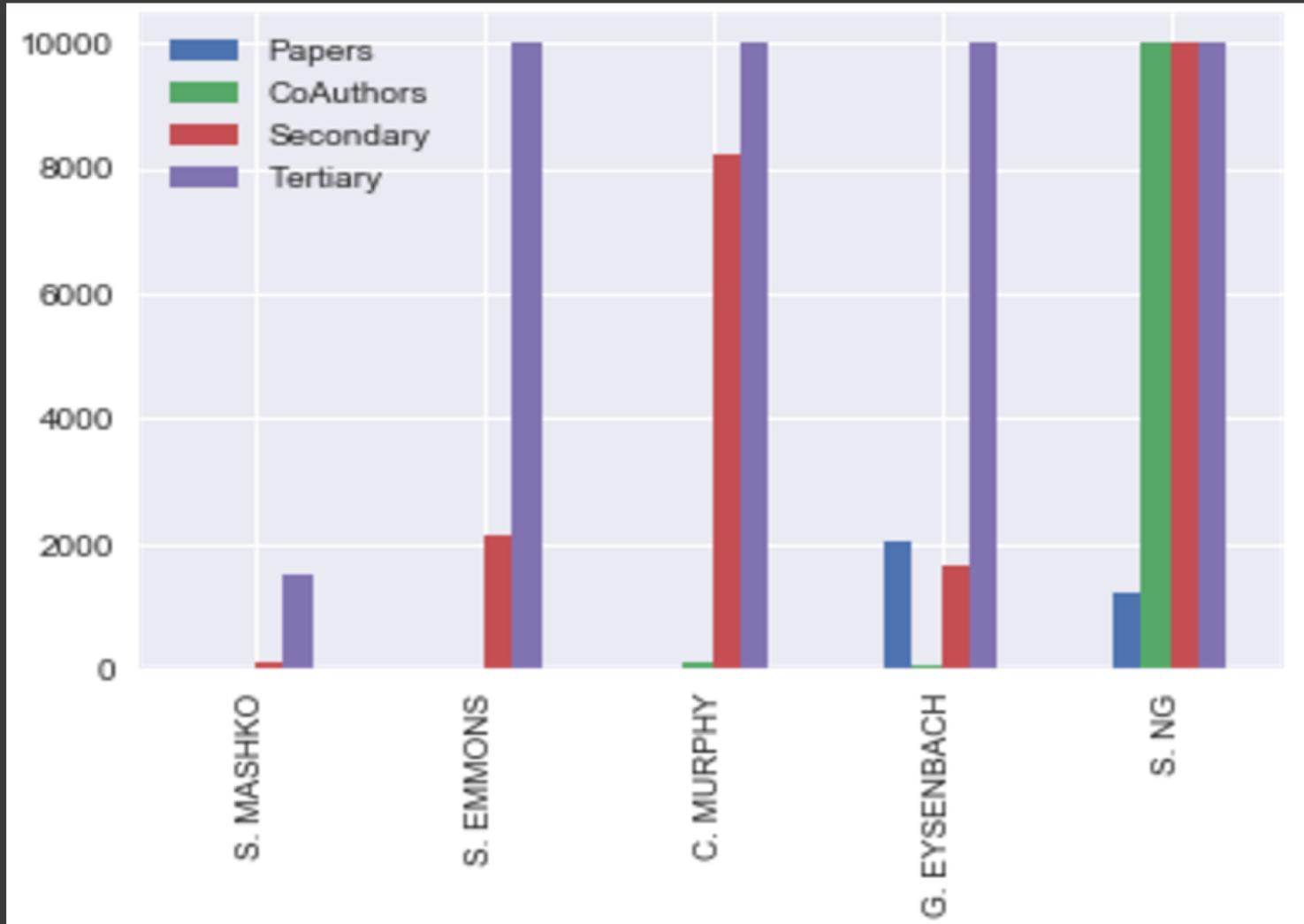
SCOTT EMMONS'  
2ND DEGREE  
COLLABORATION  
NETWORK



COLEEN MURPHY'S  
2ND DEGREE  
COLLABORATION  
NETWORK



# INCREASING NODE COUNT WITH DISTANCE FROM AUTHOR

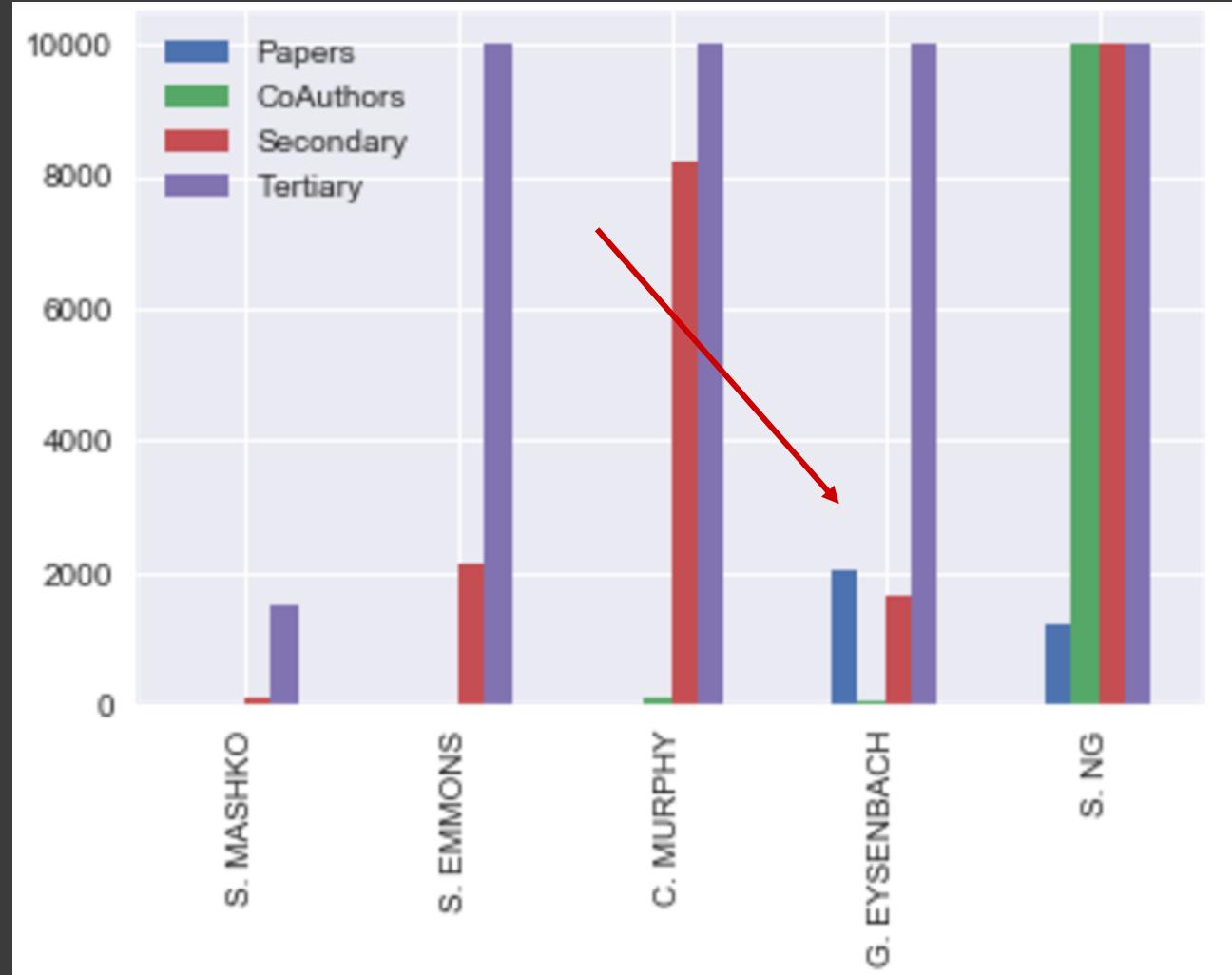


PAPERS < COAUTHORS <  
SECONDARY < TERTIARY



PAPERS <  
COAUTHORS <  
SECONDARY <  
TERTIARY

# 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.



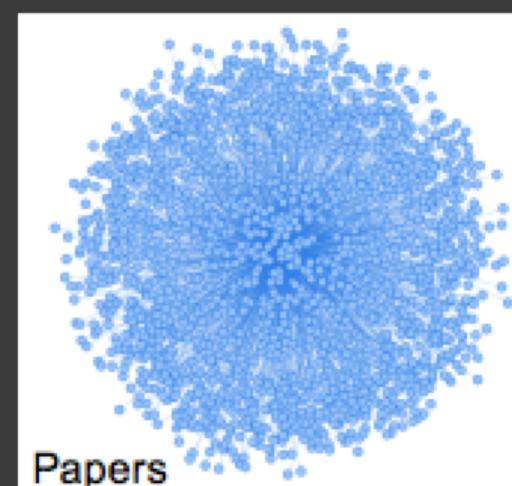
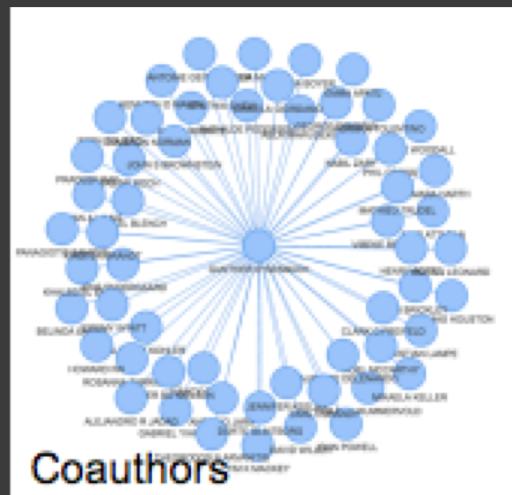
# FEW CO-AUTHORS AND MANY PAPERS SUSPICIOUS PATTERN

BUT HE IS ONLY CITED FOR 120 PAPERS AND 40 BOOK CHAPTERS.



SO WHAT ARE THOSE OTHER 1900 LINKS ?

PROBABLY EDITING JOBS (FALSE POSITIVES)



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/JANUSGRAPH



- APACHE PROJECT (TITAN)/LINUX FOUNDATION (JANUSGRAPH)
- 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>



# GRAKN.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://grakn.ai/>



# Thank You!

• Denis Vrdoljak, MIDS  
Managing Director, BDSG  
Marketing Analyst/Data Scientist, Cisco

[• denis@bds.group](mailto:denis@bds.group)

[• dvrdolja@cisco.com](mailto:dvrdolja@cisco.com)

Gunnar Kleemann, PhD, MIDS  
Senior Data Scientist, BDSG  
Data Science Instructor, UC Berkeley

[• gunnar@bds.group](mailto:gunnar@bds.group)

[• gunnarkl@berkeley.edu](mailto:gunnarkl@berkeley.edu)

