# Beyond Word2Vec:
# Using embeddings to chart
# out the ebb and flow of tech skills

**Maryam Jahanshahi Ph.D.**
Research Scientist
TapRecruit.co

tap'Recruit.co

http://bit.ly/aiconf-2019

Job ⌄ | ↻ Sync | Similar Jobs ⌄ | **Open** | Large Candidate Pool | 📊 Applicants: 202 ⌄ | 3850 Characters | Notify ⌄ | Last edit: **System** ⌄

📊 This job has not attracted enough qualified applicants.

| Days Live: | 50 | Applied: | 202 |
| Screening: | 14 | Assessment: | 1 |
| Offer: | 0 | Hired: | 0 |

Job Report →

**28** **Job will perform poorly**

This job scores **lower than 95%** of **Junior Accounting** jobs in **Los Angeles, CA**

- Add preferred qualifications
- Add more "you" statements
- Perks included
- Equal opportunity statement is included

Neutral — Gendered

# Senior Finance

TapRecruit  -  Los Angeles

TapRecruit is looking for a smart, detail-oriented person to serve as a senior financial analyst. This person will be responsible for supporting the company's FP&A requirements. Responsibilities will include working on TapRecruit Entertainment Group's FP&A model, supporting analysis for long-term plan or options, tracking key business operational metrics and producing monthly financial/operation... itional FP&A needs, this role will require strong organizational skills to help manage the ... uide discussions with senior managers across the department and evaluate/implem... ojects for top management. This is a dynamic role that serves the finance de... report to a Senior Manager of Finance and will routinely interface with TapRecruit's top ma...

Language that emphasizes an "intense" or "confusing" environment is known to deter qualified candidates.

🗑 Delete

This is an ideal position for an individual who has gained stron... investment bank or accounting firm and now seeks to apply those skills to a fast-growing entrepreneurial company. Strong quantitative and excel financial modeling skills are a must. The ideal candidate must be comfortable in a dynamic start-up environment, will bring energy and passion to everything he/she does, and will not be afraid to roll up his/her sleeves to tackle challenging analytical assignments.

This job is full-time, based in Los Angeles. We offer competitive compensation and stock option program.

# Skills and qualifications matter in job descriptions

Same title,
Different job

**Finance Manager**
**Kraft Foods**

Junior (3 Years)

No Managerial Experience

**Finance Manager**
**Roche**

Senior (6-8 Years)

Division Level Controller

Strategic Finance Role

MBA / CPA

✓ **Same Title**

✗ Required Experience
✗ Required Responsibility
✗ Preferred Skill
✗ Required Education

tapRecruit.co

# Research at TapRecruit

Helping companies make fairer and more efficient recruiting decisions

## NLP and Data Science:

- What are distinguishing characteristics of successful career documents?

- What skills are increasingly important for different industries?

## Decision Science:

- How do candidates make decisions about which jobs to apply to?

- How do hiring teams make decisions about candidate qualifications?

tapRecruit.co

# How have tech skills changed over time?

# Strategies to identify changes among corpora

Traditional approaches do not capture syntactic and semantic shifts

| MBA | SQL |
|-----|-----|
| PhD | Tableau |
| Python | PowerBI |

| 1880 | 1920 | 1960 | 2000 |
|------|------|------|------|
| **force** | **atom** | **radiat** | **state** |
| **energy** | **theory** | **energy** | **energy** |
| **motion** | **electron** | **electron** | **electron** |
| **differ** | **energy** | **measure** | **magnet** |
| **light** | **measure** | **ray** | **field** |

Matter                                    Quantum

Electron

**Manual Feature Extraction**

Require selection of key attributes, therefore difficult to discover new attributes

**Dynamic Topic Models**

Require experimentation with topic number

Adapted from Blei and Lafferty, ICML 2006.

tapRecruit.co

# Embeddings use context to extract meaning

Window sizes capture semantic similarity vs semantic relatedness

Statistical modeling through software (e.g. SPSS) or programming language (e.g. **Python**)

**Context** — **Word**

Experience in **Python**, Java or other object-oriented programming languages

**Context** **Word** **Context**

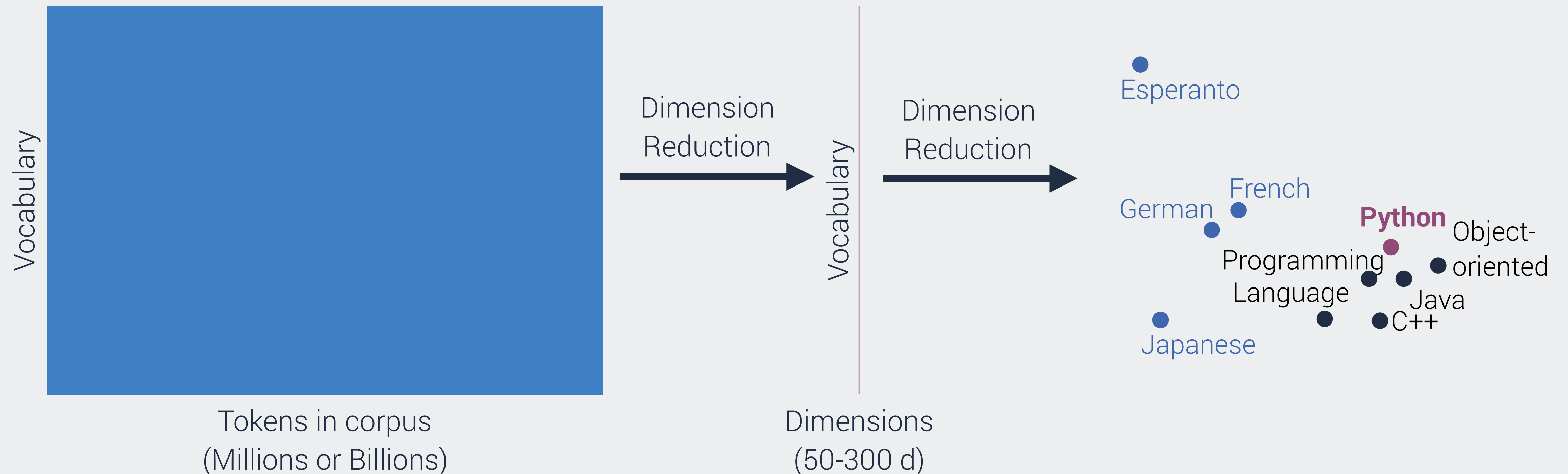Proficiency programming in **Python**, Java or C++.

**Context** **Word** **Context**
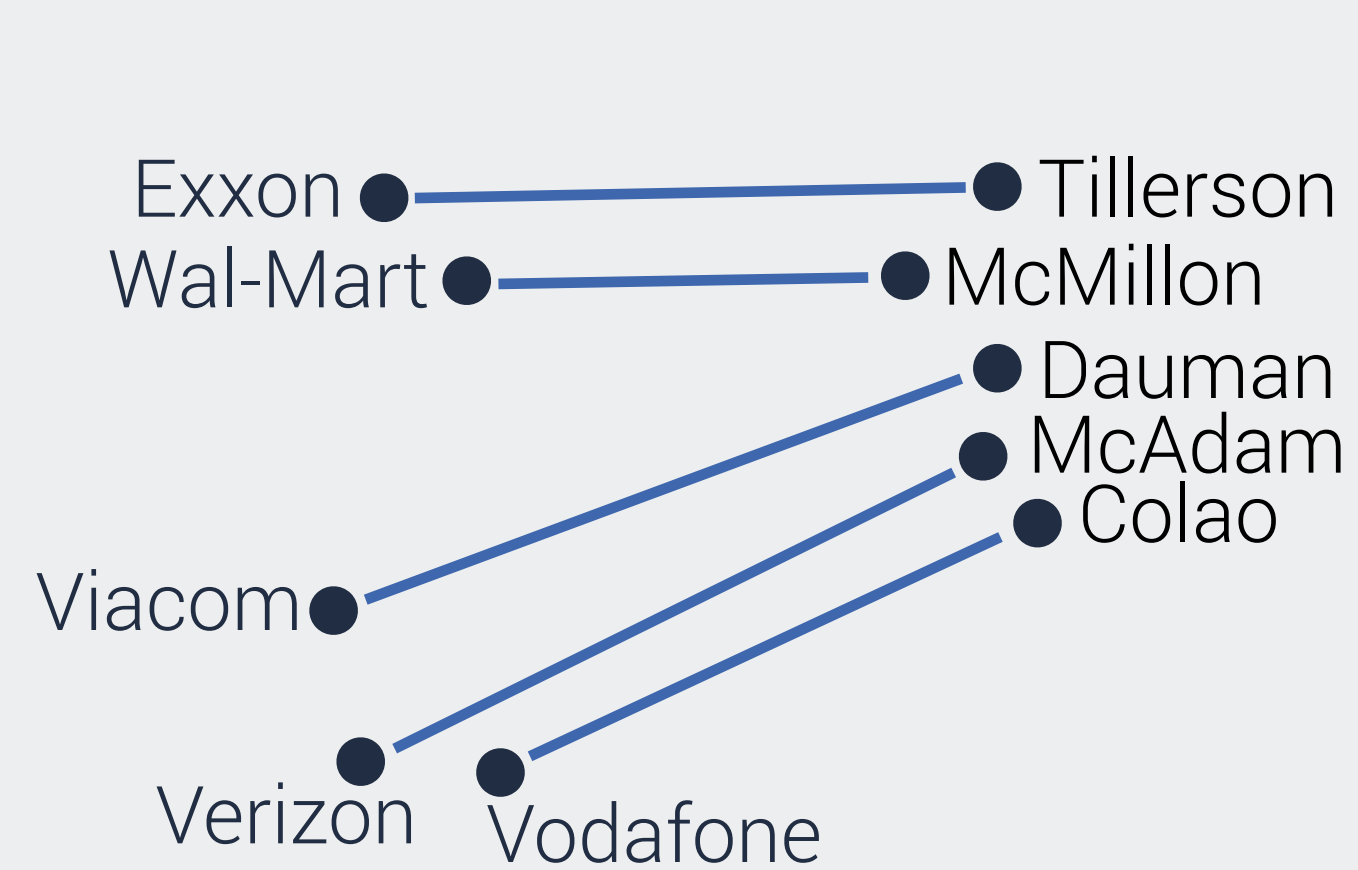
tapRecruit.co

# A simplified representation of word vectors

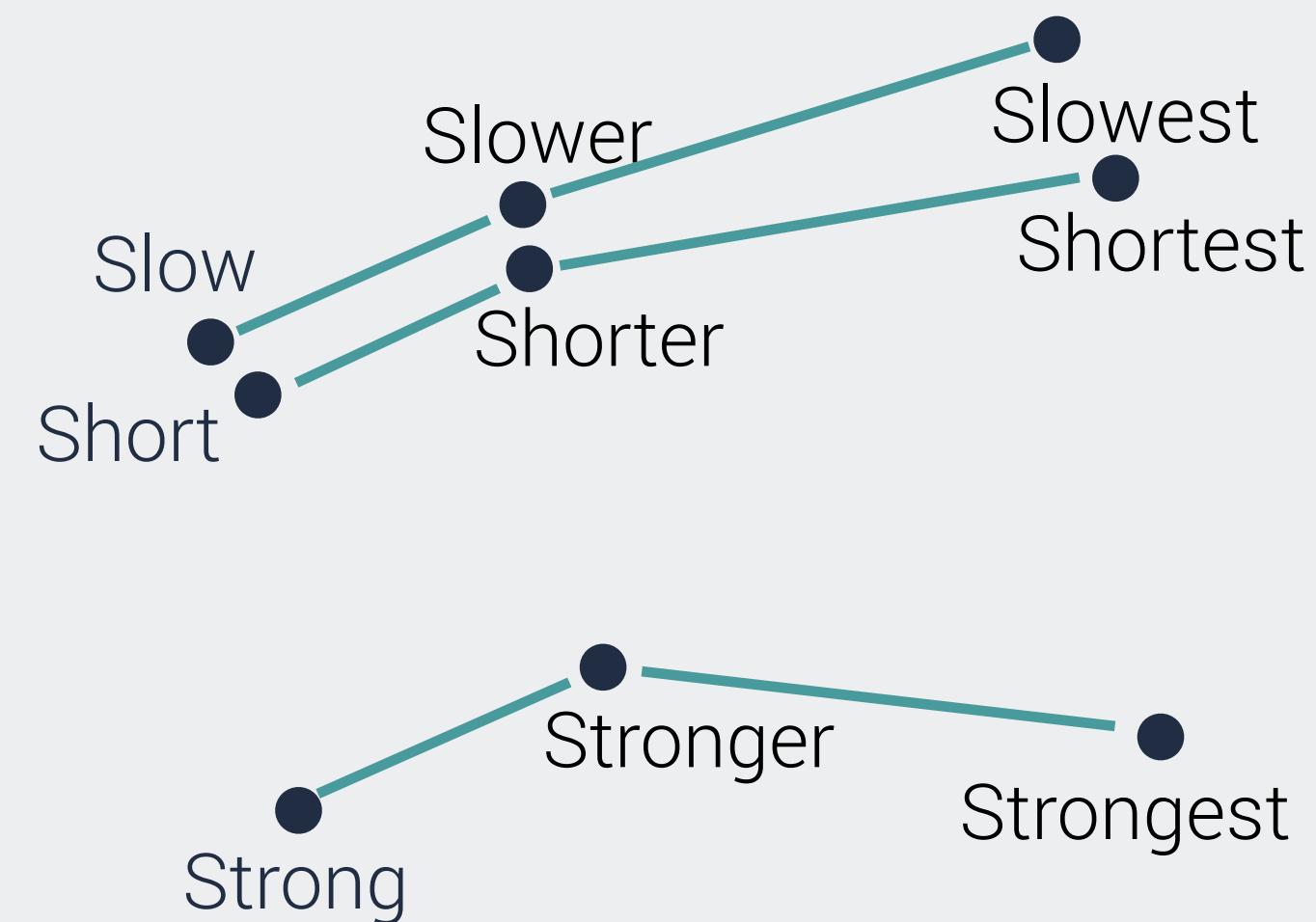Dimension reduction is key to all types of embeddings models

# Embeddings capture entity relationships

Dimensionality enables comparison between word pairs along many axes
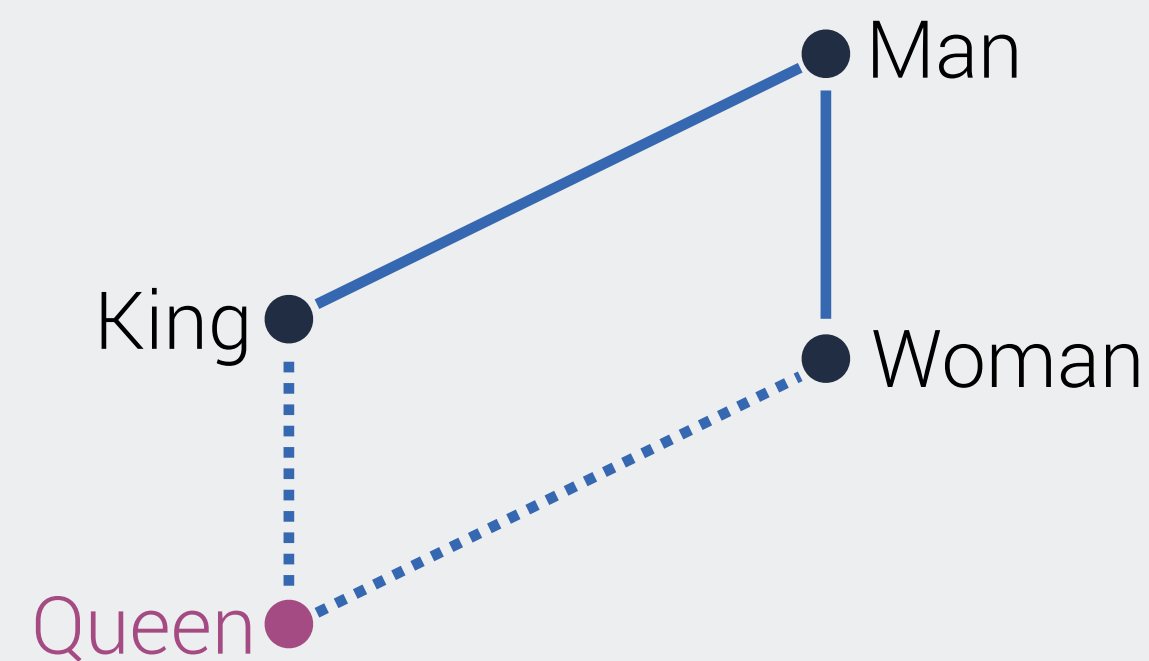


**Hierarchies**
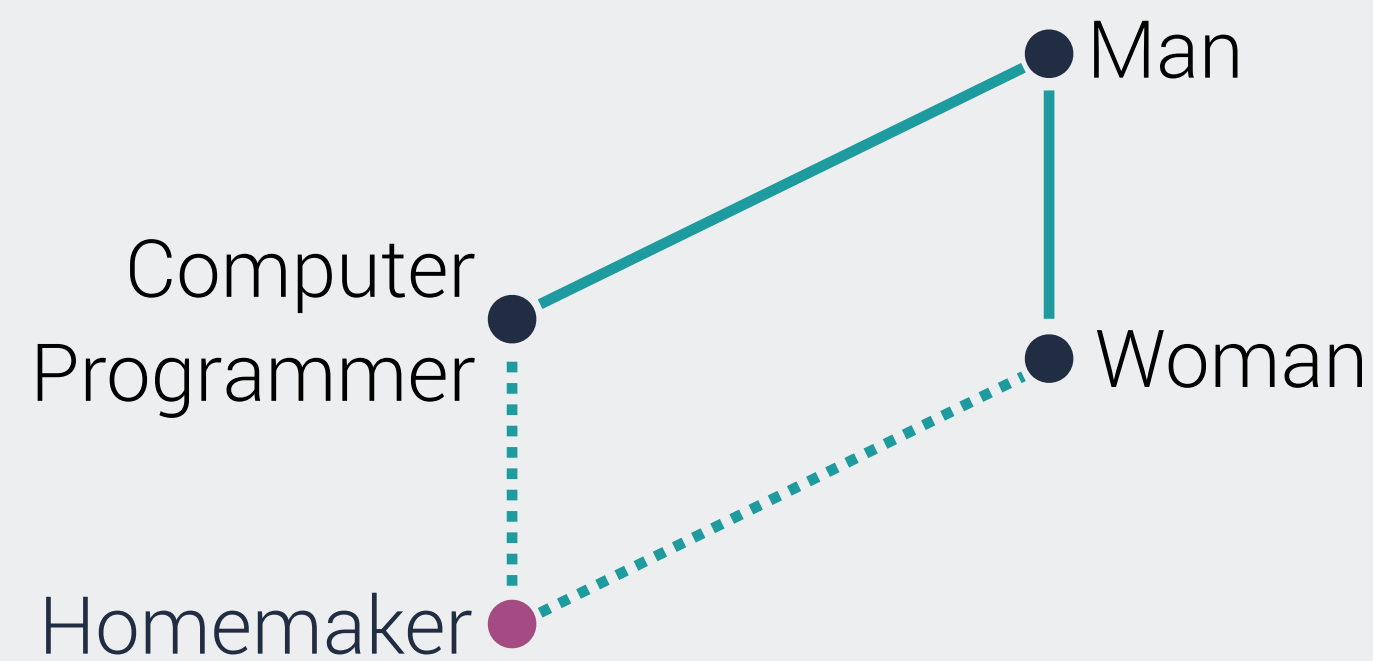
**Comparatives and Superlatives**

**Man :: King as Woman :: ?**

Adapted from Stanford NLP GLoVE Project

tapRecruit.co

# Embeddings reflect cultural bias in corpora

## High dimensionality enables some bias reduction



**Man :: King as Woman :: ?**

**Man :: Programmer as Woman :: ?**

Adapted from Bolukbasi et al., arXiv: 1607:06520.

tapRecruit.co

# Pretrained embeddings facilitate fast prototyping

Embeddings training should match corpus that is being tested on

| Corpus Generation | Corpus Tokens | Twitter 27 B | Common Crawl 42-840 B | GoogleNews 100 B | Wikipedia 6 B |
|---|---|---|---|---|---|
| Corpus Processing | Vocabulary Size | 1.2 M | 1.9-2.2 M | 3 M | 400 k |
| Language Model Generation | Algorithm Vector Length | GLoVE 25 - 200 d | GLoVE 300 d | word2vec 300 d | GLoVE 50 - 300 d |
| Language Model Tuning | | | | | |
| Final Application | | | | | |

tapRecruit.co

# Problems with pretrained embedding models

| | |
|---|---|
| **Casing** | Abbreviations vs Words<br>e.g. IT vs it |
| **Out of Vocabulary Words** | Domain Specific Words & Acronyms |
| **Polysemy** | Words with multiple meanings<br>e.g. drive (a car) vs drive (results)<br>e.g. Chef (the job) vs Chef (the language) |
| **Multi-word Expressions** | Phrases that have new meanings<br>e.g. Front-end vs front + end |

**tap**Recruit.co

# Custom language models tools

Modularized for different data and modeling requirements

**spaCy**  **OpenNLP™**

**CoreNLP**   SyntaxNet

**gensim**  **PYTORCH**

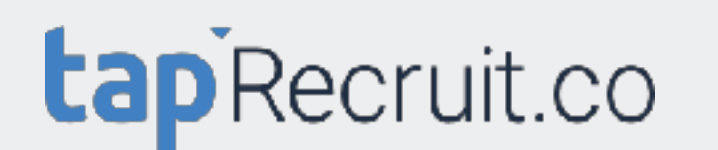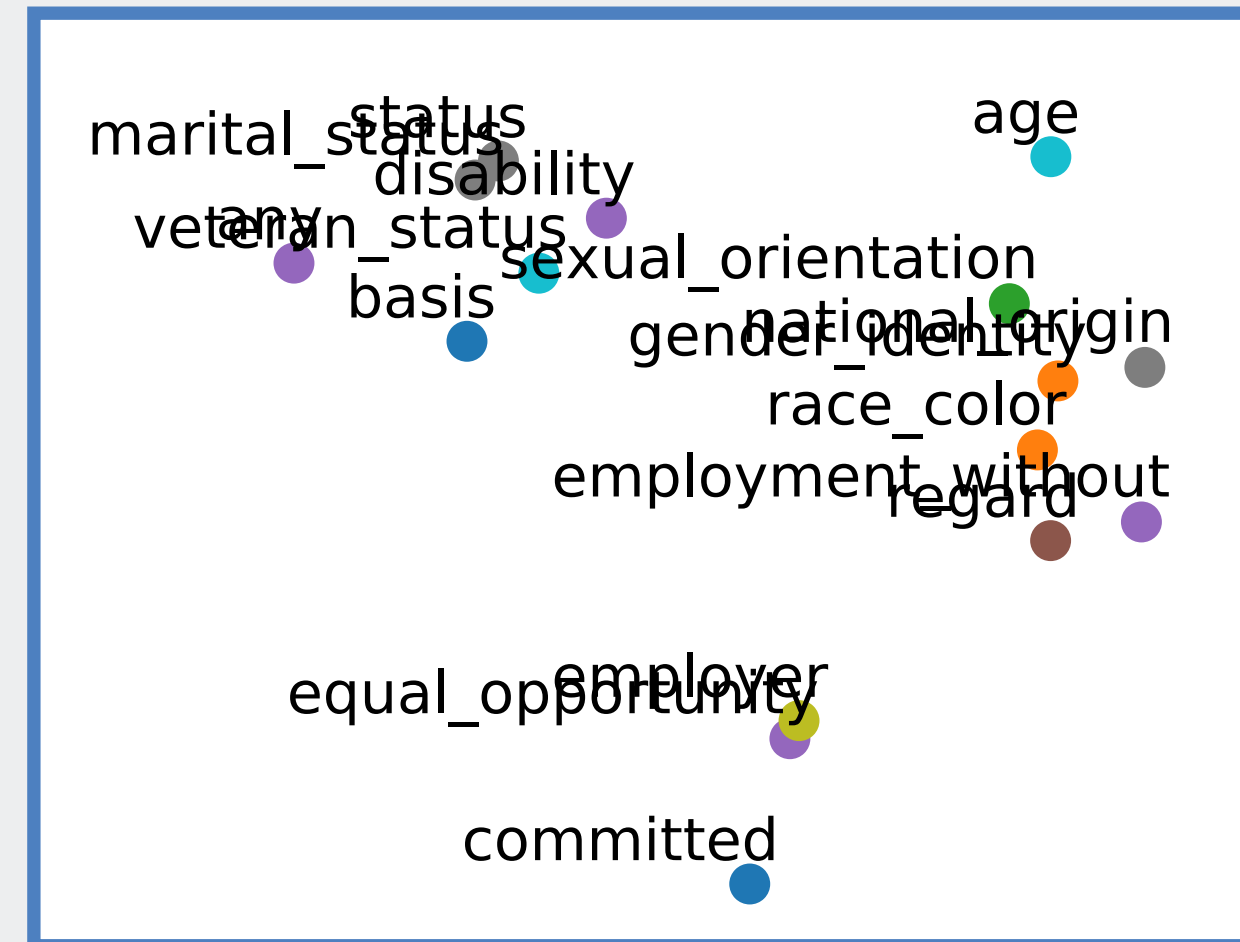**TensorFlow**  Amazon SageMaker

**Corpus Processing**
Tokenization, POS tagging, Sentence
Segmentation, Dependency Parsing

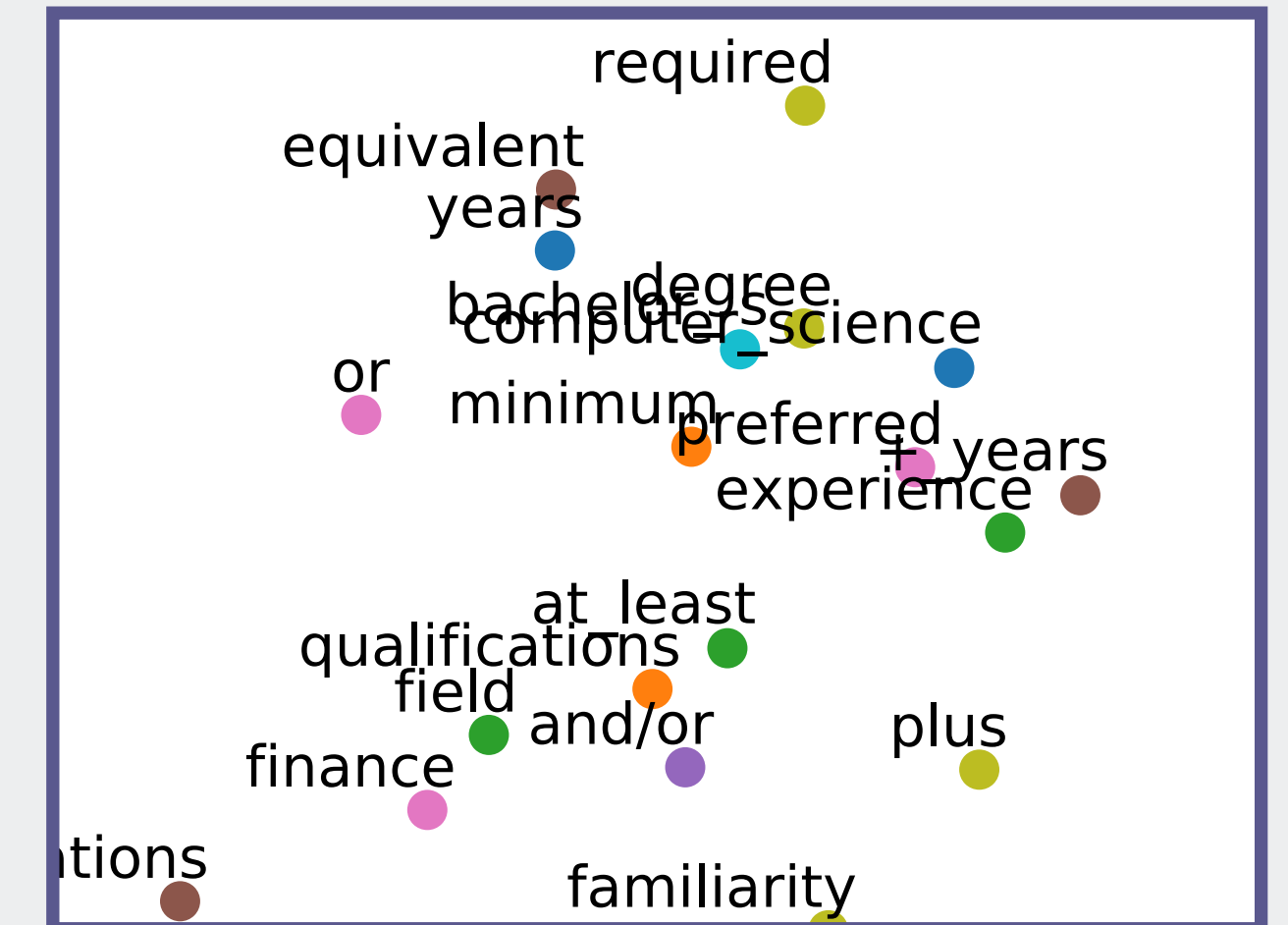**Language Modeling**
Different word embedding models
(GLoVE, word2vec, fastText)

tapRecruit.co

# Career language embedding model

## Identified equal opportunity and perks language

# Career language embedding model
## Identified 'soft' skills and language around experience



tapRecruit.co

I've got 300 dimensions…
but time ain't one

# Two approaches to connect embeddings

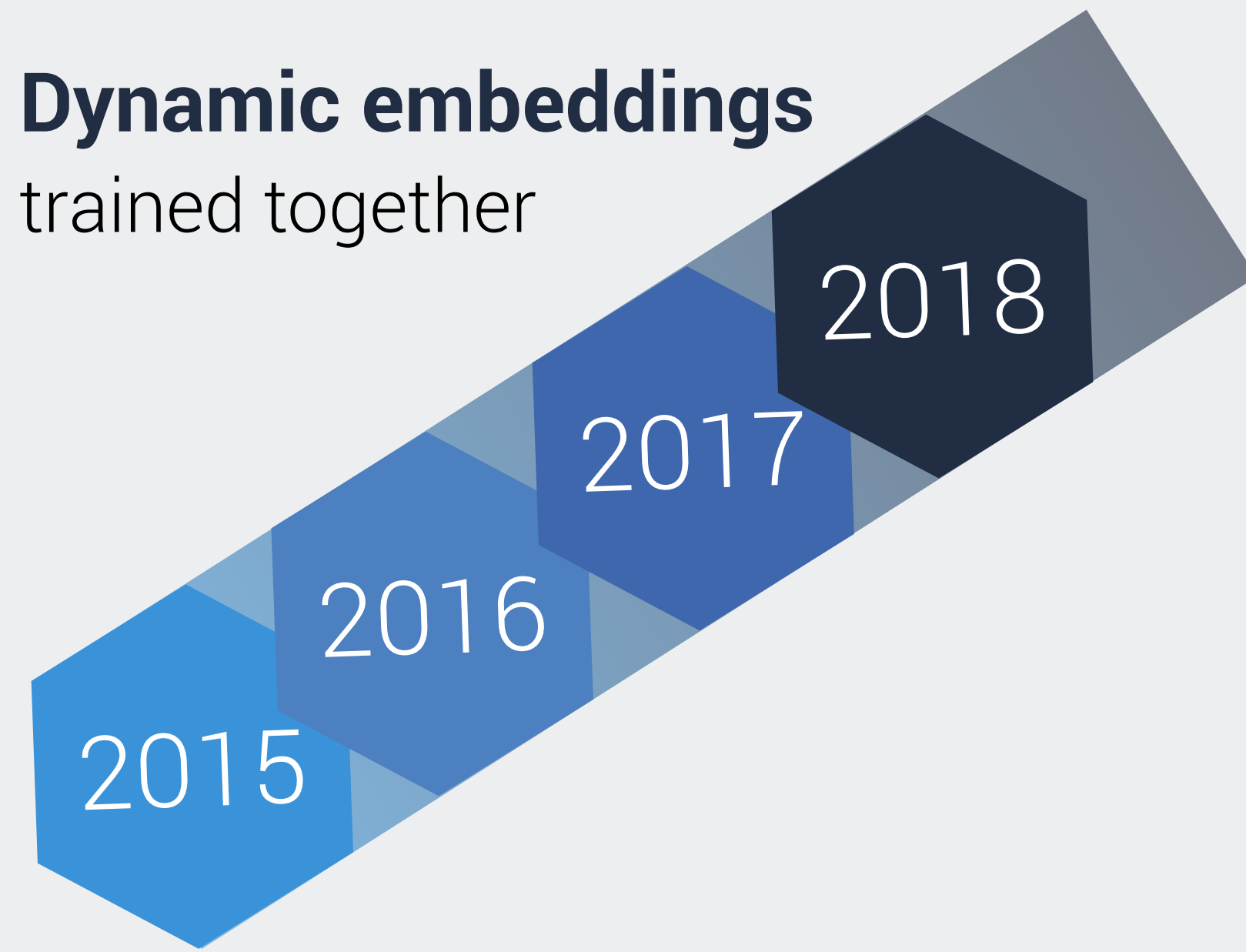**Static embeddings**
stitched together

2018

2017

2016

2015

**Data hungry**
**Requires alignment**

Kim, Chiu, Kaneki, Hedge and Petrov, arXiv: 1405:3515.
Kulkarni, Al-Rfou, Perozzi and Skiena, arXiv: 1411:3315.

**Dynamic embeddings**
trained together

2018

2017

2016

2015

**Data efficient**
**Does not require alignment**

Balmer and Mandt, arXiv: 1702:08359
Yao, Sun, Ding, Rao and Xiong, arXiv: 1703:00607
**Rudolph and Blei, arXiv: 1703:08052**

tapRecruit.co

# Dynamic Bernoulli embeddings

Outputs facilitate quick analysis of trends

## Absolute drift

Identifies top words whose usage
changes over time course

| words with largest drift (Senate) | | | |
|---|---|---|---|
| IRAQ | 3.09 | coin | 2.39 |
| tax cuts | 2.84 | social security | 2.38 |
| health care | 2.62 | FINE | 2.38 |
| energy | 2.55 | signal | 2.38 |
| medicare | 2.55 | program | 2.36 |
| DISCIPLINE | 2.44 | moves | 2.35 |
| text | 2.41 | credit | 2.34 |
| VALUES | 2.40 | UNEMPLOYMENT | 2.34 |

## Embedding neighborhoods

Extract semantic changes by nearest
neighbors of drifting words

| UNEMPLOYMENT | | |
|---|---|---|
| **1858** | **1940** | **2000** |
| unemployment | unemployment | unemployment |
| unemployed | unemployed | jobless |
| depression | depression | rate |
| acute | alleviating | depression |
| deplorable | destitution | forecasts |
| alleviating | acute | crate |
| destitution | reemployment | upward |
| urban | deplorable | lag |
| employment | employment | economists |
| distressing | distress | predict |

**Repository Link:** http://bit.ly/dyn_bern_emb

# Experiments with dynamic embeddings

| | Small Corpus |
|---|---|
| **Job Types** | All |
| **Time Slices** | 3 (2016-2018) |
| **Number of Documents** | **50 k** |
| **Vocabulary Size** | 10 k |
| **Data Preprocessing** | Basic |
| **Embedding Dimensions** | 100 d |

tapRecruit.co
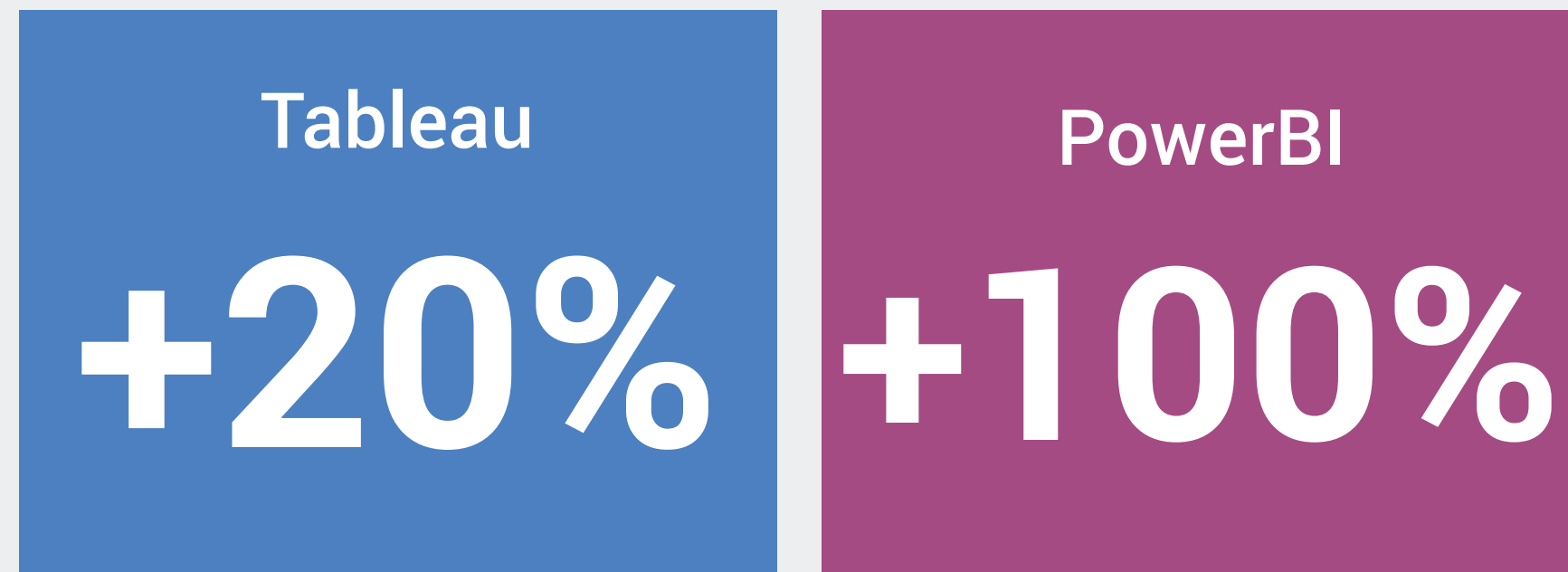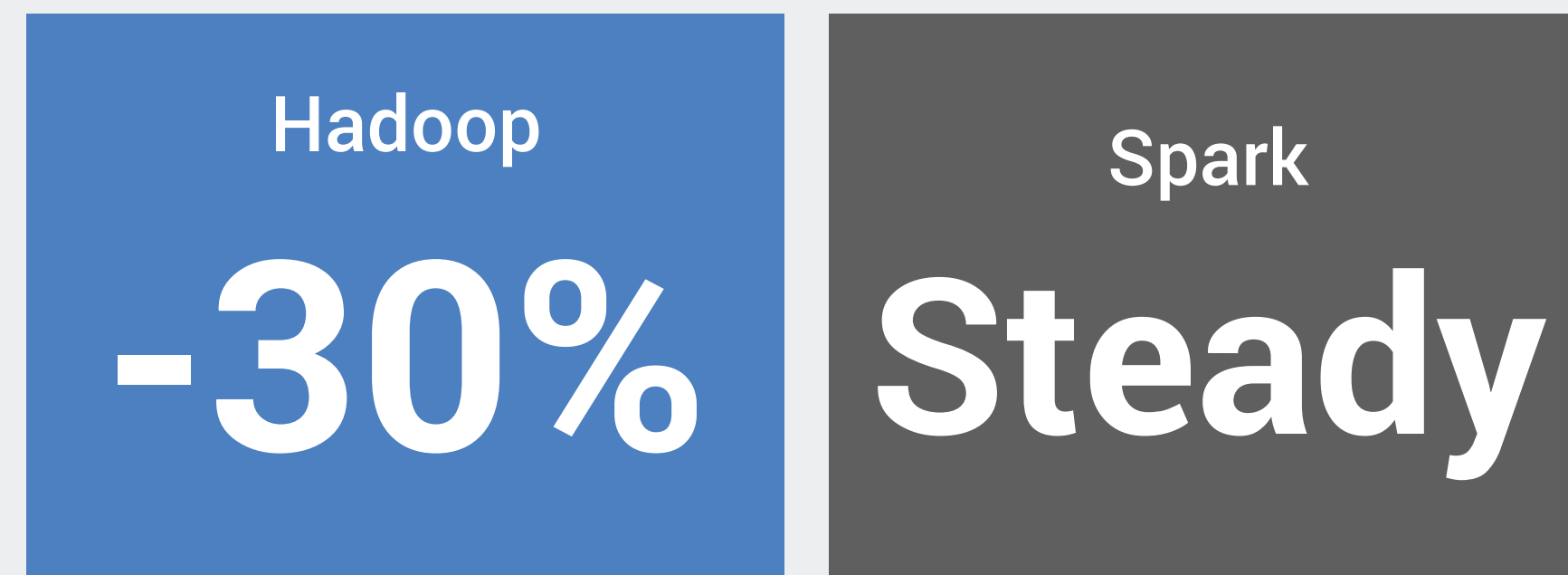
# Small corpus identified skill demands

Data Viz is up and Hadoop (but not Spark) is down

**Demand for Data Visualization tools is up**

| Tableau | PowerBI |
|---------|---------|
| **+20%** | **+100%** |

**Demand for Hadoop is down in DS and ML roles**

| Hadoop | Spark |
|--------|-------|
| **-30%** | **Steady** |

Blue boxes indicate phrases identified from top drifting words analysis. Grey and pink boxes indicate 'control' skills.

tapRecruit.co

# Battle of the Languages

Difference between supply vs demand of scripting languages

## Demand for Perl is down

| Perl | Python |
|------|--------|
| **-40%** | **Steady** |



stack**overflow**

Developer Survey Results
2019

- Python, the fastest-growing major programming language, has risen in the ranks of programming languages in our survey yet again, edging out Java this year and standing as the second most loved language (behind Rust).
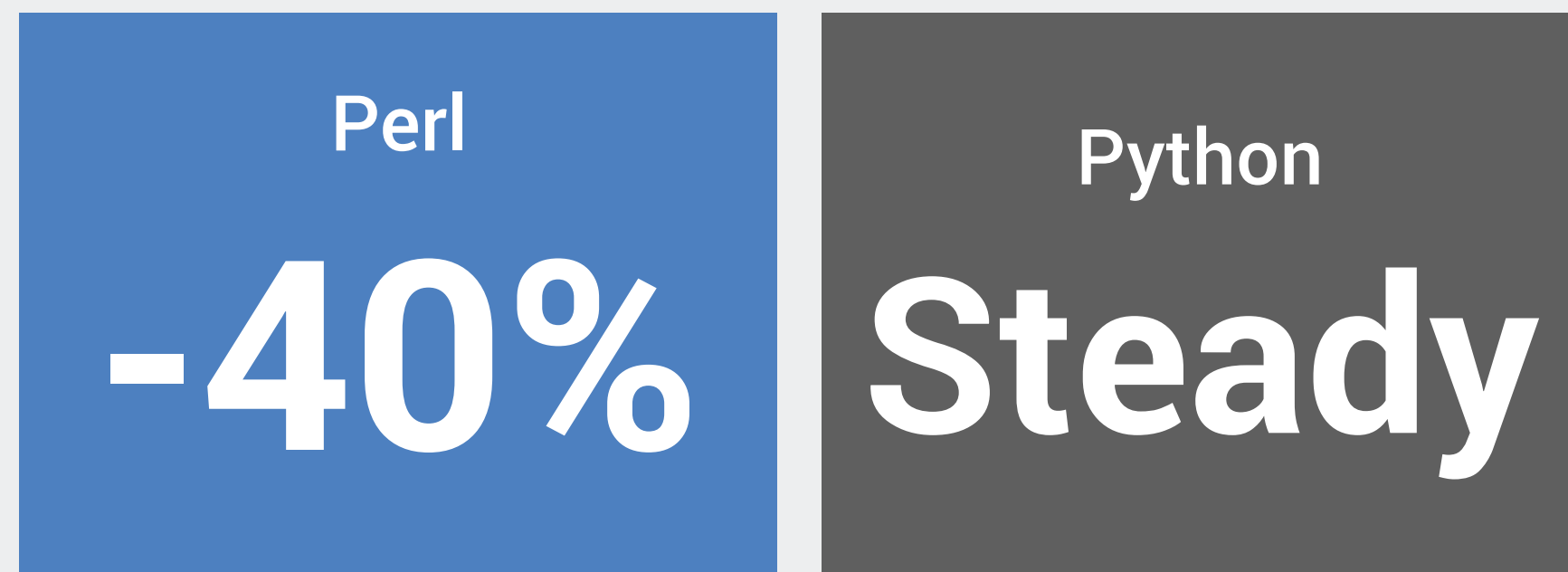
Blue boxes indicate phrases identified from top drifting words analysis. Grey and pink boxes indicate 'control' skills. **tap**Recruit.co

# Battle of the Languages
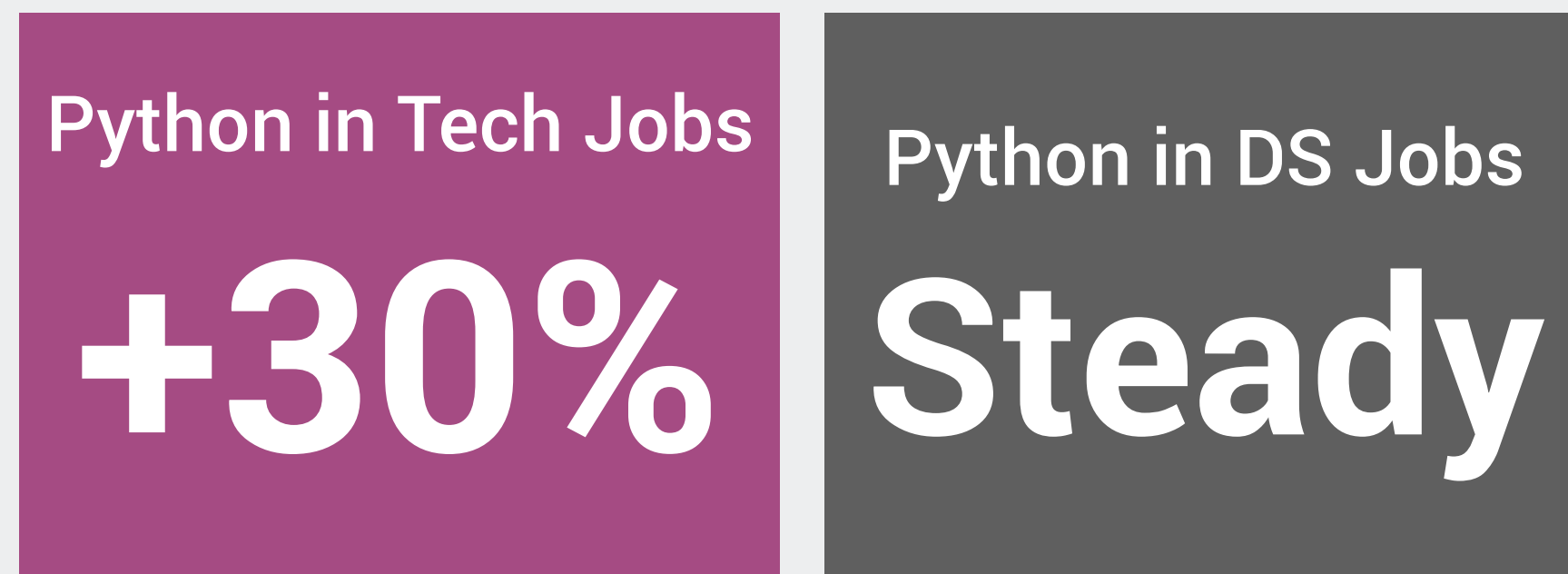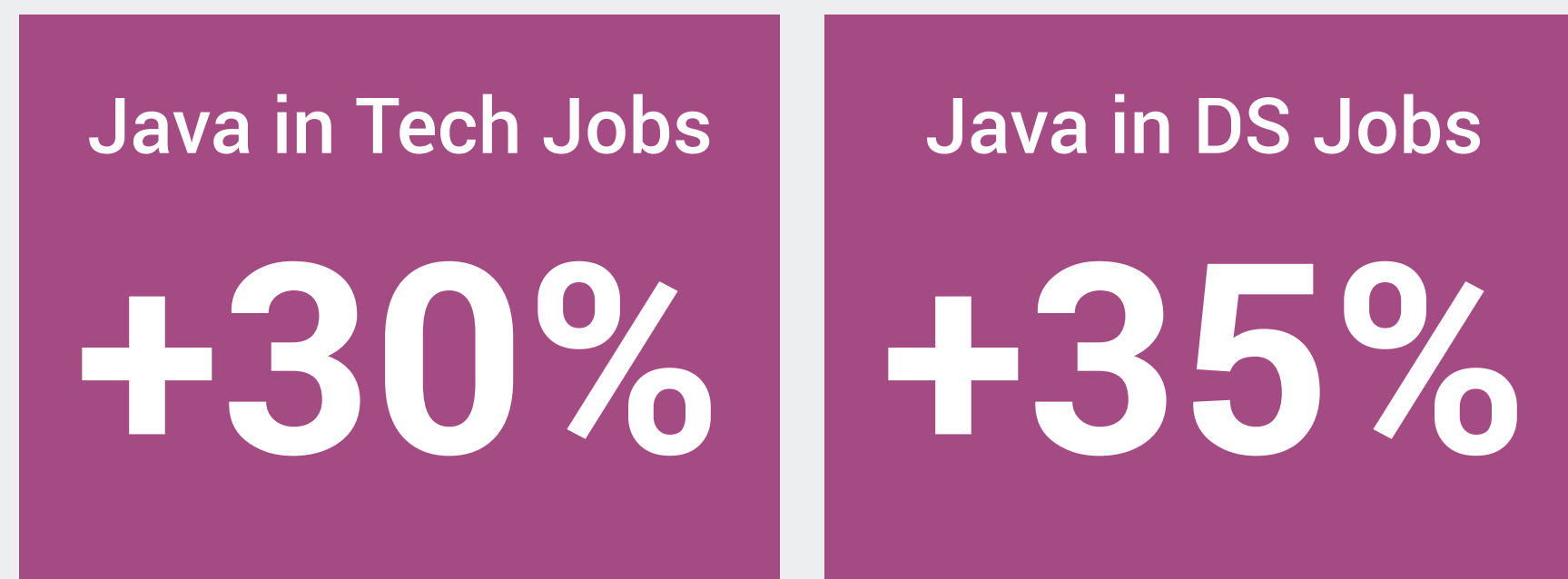
Difference between supply vs demand of scripting languages

## Demand for Python up in Tech roles

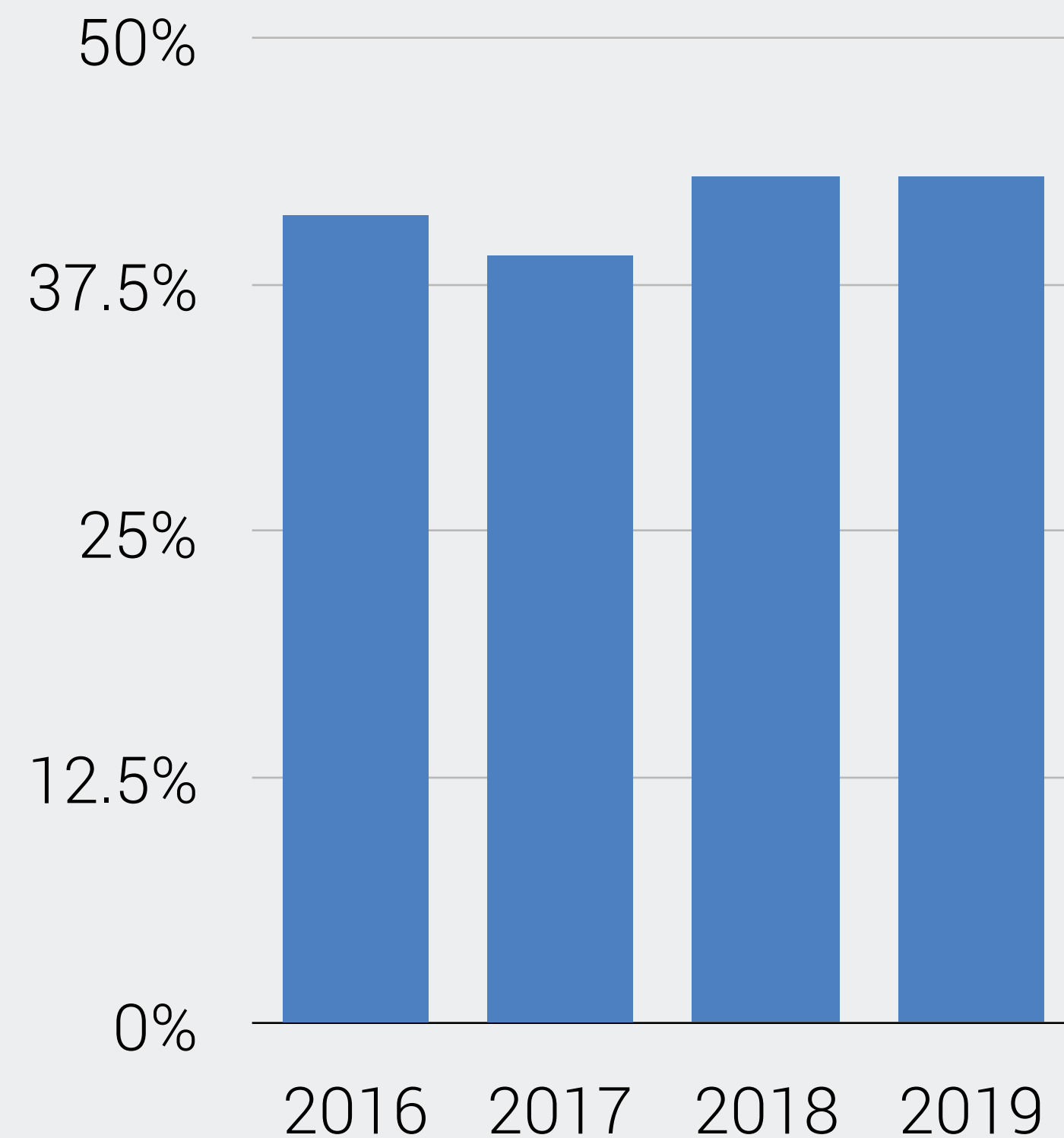| Python in Tech Jobs | Python in DS Jobs |
|---|---|
| **+30%** | **Steady** |

## Demand for Java is up

| Java in Tech Jobs | Java in DS Jobs |
|---|---|
| **+30%** | **+35%** |



stack overflow

Developer Survey Results
2019

- Python, the fastest-growing major programming language, has risen in the ranks of programming languages in our survey yet again, edging out Java this year and standing as the second most loved language (behind Rust).

Blue boxes indicate phrases identified from top drifting words analysis. Grey and pink boxes indicate 'control' skills. tapRecruit.co

# Experiments with dynamic embeddings

| | Small Corpus | Large Corpus |
|---|---|---|
| Job Types | All | All |
| Time Slices | 3 (2016-2018) | 3 (2016-2018) |
| Number of Documents | **50 k** | **500 k** |
| Vocabulary Size | 10 k | 10 k |
| Data Preprocessing | Basic | Basic |
| Embedding Dimensions | 100 d | 100 d |

# Beyond word2vec

- Flavors of static word embeddings: The Corpus Issue

- Considerations for developing custom embedding models

- Dynamic Bernoulli embeddings are robust with small datasets

# How have tech and data science skills changed?

- Demand for MBAs and PhDs is falling

- Core Skills: DataViz & Scripting Languages

- Commodification of distributed systems impacts demand for Hadoop

- Demand for SQL in a variety of core business functions

# Thank you AI Conference!

**Maryam Jahanshahi Ph.D.**

Research Scientist

🐦 @mjahanshahi

in maryam-j

**tap**'Recruit.co

http://bit.ly/aiconf-2019