# Introduction to NLP

*Searching for Meaning in Text*

Dr. Tony Russell-Rose FBCS CITP CEng

uxlabs

Goldsmiths
UNIVERSITY OF LONDON

2Dsearch

# Introductions

**Tony Russell-Rose, PhD**

Senior Lecturer, Goldsmiths    Director, UXLabs

**uxlabs**

- Led R&D / technology innovation at Microsoft, Canon, Reuters, BT Labs, HP Labs, Oracle

- Visiting Professor of Cognitive Computing & AI, Essex University

- PhD (natural language interfaces), MSc in human-computer interaction, first degree in engineering (human factors)

- 5 patents, 80+ scientific and technical papers on AI/NLP, information retrieval and UX

- Founder, 2Dsearch: next generation advanced search

- Director, UXLabs: UX research & design consultancy

- Honorary Visiting Fellow at City University Centre for Interactive Systems Research

- Founder of Search Solutions conference series. Served as Vice-chair of BCS IRSG, Chair of CIEHF HCI Group

- Blog: http://isquared.wordpress.com

**Industry Expertise**

- Media & publishing
- Healthcare
- Business intelligence
- eCommerce

**Domain Expertise**

- UX research & design
- Information architecture
- Search & information retrieval
- AI + machine learning

**2Dsearch**

# NLP – a solved problem?

- As humans we do it effortlessly … don't we?

- DRUNK GETS NINE YEARS IN VIOLIN CASE
- PROSTITUTES APPEAL TO POPE
- STOLEN PAINTING FOUND BY TREE
- RED TAPE HOLDS UP NEW BRIDGE
- DEER KILL 300,000
- RESIDENTS CAN DROP OFF TREES
- INCLUDE CHILDREN WHEN BAKING COOKIES
- MINERS REFUSE TO WORK AFTER DEATH

# Problems with Text

- Polysemy
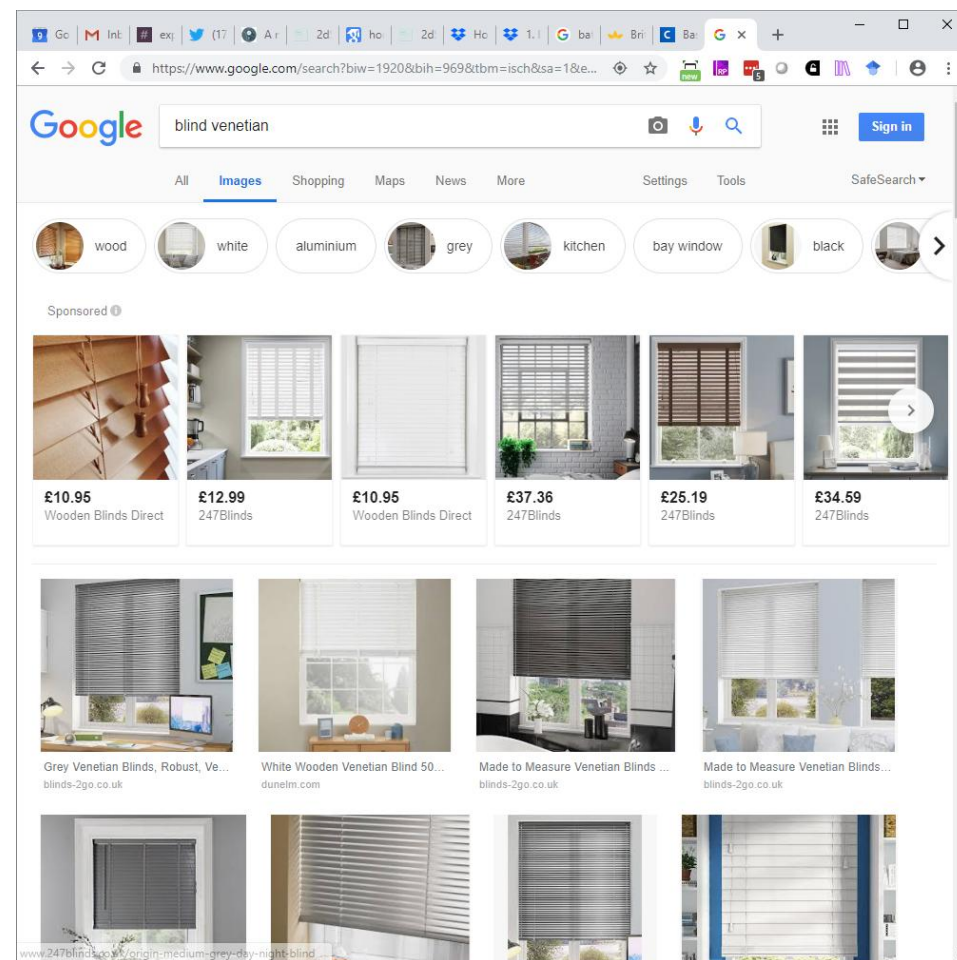  - One word maps to many concepts
  - e.g. *bat*
- Synonymy
  - One concept maps to many words

# Problems with Text

- ## Word order

  - ## *Venetian blind vs. Blind venetian*



Intro to Natural Language Processing

# Problems with Text

- ## Language is generative
  - `Starbucks coffee is the best`
  - `The place I like most when I need to feed my caffeine addiction is the company from Seattle with branches everywhere`
- ## Many different ways to express a given idea
  - Synonymy, paraphrase, metaphor, etc.

# Problems with Text

- Language is changing
  - `I want to buy a mobile`
- Ill-formed input
  - `"accomodation office"`
- Co-ordination, negation, etc.
  - `This is not a talk about neuro-linguistic programming`
- Multi-linguality
  - `Claudia Schiffer is on the cover of Elle`
- Sarcasm, irony, slang, jargon, etc.
  - `That was a wicked lecture`
  - `Yep – the coffee break was the best part`

Dan Jurafsky

# Why else is natural language understanding difficult?

## non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either♥

## segmentation issues

the | New | York-New | Haven | Railroad
the | New | York-New | Haven | Railroad

## idioms

dark horse
get cold feet
lose face
throw in the towel

## neologisms

unfriend
Retweet
bromance

## world knowledge

Mary and Sue are sisters.
Mary and Sue are mothers.

## tricky entity names

Where is *A Bug's Life* playing …
*Let It Be* was recorded …
… a mutation on the *for* gene …

But that's what makes it fun!

# NLP Fundamentals (word level)

- ## Language is AMBIGUOUS
  - To determine structure, we must resolve ambiguity!

- ## Lexical analysis (tokenisation)
  - `The cat sat on the mat`
  - `I can't tokenise this sentence`

- ## Stop word removal
  - No definitive list
    - `The Who, The The, Take That…`
    - `To be or not to be`

# NLP Fundamentals (word level)

- ## Stemming
  - `fishing, fished, fish, fisher -> fish`
- ## Lemmatization
  - ### Linguistically principled analysis
  - `Passing -> pass + ING`
  - `Were -> be + PAST`
  - `Delegate = de-leg-ate (?)`
  - `Ratify = rat-ify (?)`

- ## Morphology (prefixes, suffixes, etc.)
  - `Gebäudereinigungsfirmenangestellter -> Gebäude + Reinigung + Firma + Angestellter (building + cleaning + company + employee)`

# NLP Fundamentals (sentence level)

- ## Syntax (part of speech tagging)
  - `book -> NOUN, VERB`
  - `that -> DETERMINER`
  - `flight -> NOUN`
  - `Book that flight -> VB DT NN`

- ## Ambiguity problem
  - `Time flies like an arrow -> ?`
  - `Fruit flies like a banana -> ?`
  - `Eats shoots and leaves -> ?`

# NLP Fundamentals (sentence level)

- # Parsing (grammar)
  - **I saw a venetian blind**
  - **I saw a blind venetian**
  - **I saw the man on the hill with a telescope**
  - **Rugby is a game played by men with odd-shaped balls**

- # Sentence boundary detection
  - **Punctuation denotes the end of a sentence!**
  - **"But not always!", said Fred...**

# NLP Fundamentals (paragraph level)

- Anaphora resolution
  - "John dropped a plate. It broke."
- Anaphora resolution relies on knowledge:
  - `We gave the bananas to the monkeys because `**`they`**` were hungry.`
  - `We gave the bananas to the monkeys because `**`they`**` were ripe.`
  - `We gave the bananas to the monkeys because `**`they`**` were here.`

What do we want?

*Anaphora resolution!*

When do we want it?

*When do we want what?*

# Information Extraction

Event:  Curriculum mtg
Date:   Jan-16-2012
Start:  10:00am
End:    11:30am
Where: Gates 159

Subject: **curriculum meeting**

Date: January 15, 2012

To: Dan Jurafsky

---

Hi Dan, we've now scheduled the curriculum meeting.

It will be in Gates 159 tomorrow from 10:00-11:30.

-Chris

Create new Calendar entry

# Named Entity Recognition

- Identification of key concepts,
  - e.g. people, places, organisations, etc.
  - Also postcodes, temporal/numerical expressions, etc.

- *"Mexico has been trying to stage a recovery since the beginning of this year and it's always been getting ahead of itself in terms of fundamentals," said Matthew Hickman of Lehman Brothers in New York."*

| Persons | Organisations | Cities | Countries |
|---|---|---|---|
| Matthew Hickman | Lehman Brothers | New York | Mexico |

# NLP Applications

- Human-computer interfaces (chatbots etc.)
- Text classification
- Text summarisation
- Machine translation
- Speech recognition & synthesis
- Natural language generation
- Text mining
- Question answering
- Sentiment Analysis

# Early NLP Systems

- ## ELIZA

  - Wiezenbaum 1966

  - Simple pattern matching

- ## SHRDLU

  - Winograd 1970

  - Natural language understanding

  - Comprehensive grammar of English

# Text Categorization

- Media, publishing, libraries…
  - Classify news stories, papers, books…
- Spam detection
- Search engines:
  - classify query intent, e.g. search Google for 'LOG313'

# Summarisation

- **Summarization**
  - single-document vs. multi-document
- **Search result snippets**
- **Word processing tools**
- **Research / analysis tools**

# Machine translation



[Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation](#), 2016.

# Dialog systems

- Chatbots
- Smart speakers
- Smartphone assistants
- Call handling systems
  - Travel
  - Hospitality
  - Banking

# Text Mining

- Analogy with *Data Mining*
  - Discover or infer new knowledge from unstructured text resources
- A <-> B and B <-> C
  - Infer A <-> C?
  - e.g. link between migraine headaches and magnesium deficiency
- Applications in life sciences, media/publishing, counter terrorism, competitive intelligence

# Question Answering: IBM's Watson

- Won Jeopardy on February 16, 2011!

WILLIAM WILKINSON'S
"AN ACCOUNT OF THE PRINCIPALITIES OF
WALLACHIA AND MOLDOVIA"
INSPIRED THIS AUTHOR'S
MOST FAMOUS NOVEL

→ Bram Stoker

# Sentiment Analysis

- Identify and extract *subjective* information
  - Predict stock market movements
- Sub-tasks:
  - Identify *polarity*, e.g. of movie reviews
    - e.g. positive, negative, or neutral
  - Identify emotional states
    - e.g. angry, sad, happy, etc.
  - Subjectivity/objectivity identification
    - E.g. "fact" from opinion
  - Feature/aspect-based
    - Differentiate between specific features or aspects of entities

Dan Jurafsky

# Language Technology

## making good progress

## mostly solved

### still really hard

### Spam detection

Let's go to Agra!   ✓

Buy V1AGRA ...   ✗

### Part-of-speech (POS) tagging

ADJ   ADJ   NOUN  VERB   ADV

Colorless  green  ideas  sleep  furiously.

### Named entity recognition (NER)

PERSON         ORG           LOC

Einstein met with UN officials in Princeton

### Sentiment analysis

Best roast chicken in San Francisco!   👍

The waiter ignored us for 20 minutes.   👎

### Coreference resolution

Carter told Mubarak he shouldn't run again.

### Word sense disambiguation (WSD)

I need new batteries for my *mouse*.

### Parsing

I can see Alcatraz from the window!

### Machine translation (MT)

第13届上海国际电影节开幕...   ⇨

The 13th Shanghai International Film Festival...

### Information extraction (IE)

You're invited to our dinner party, Friday May 27 at 8:30

Party
May 27
add

### Question answering (QA)

Q. How effective is ibuprofen in reducing fever in patients with acute febrile illness?

### Paraphrase

XYZ acquired ABC yesterday

ABC has been taken over by XYZ

### Summarization

The Dow Jones is up

The S&P500 jumped   ⇨   Economy is good

Housing prices rose

### Dialog

Where is Citizen Kane playing in SF?

Castro Theatre at 7:30. Do you want a ticket?

# Platforms & Toolkits

- Many commercial vendors

- Freely available:
  - GATE (Sheffield University)
  - Stanford CoreNLP
  - Apache OpenNLP
  - TextBlob
  - Textacy
  - Spacy
  - Nltk

- Also: RapidMiner, R, Weka
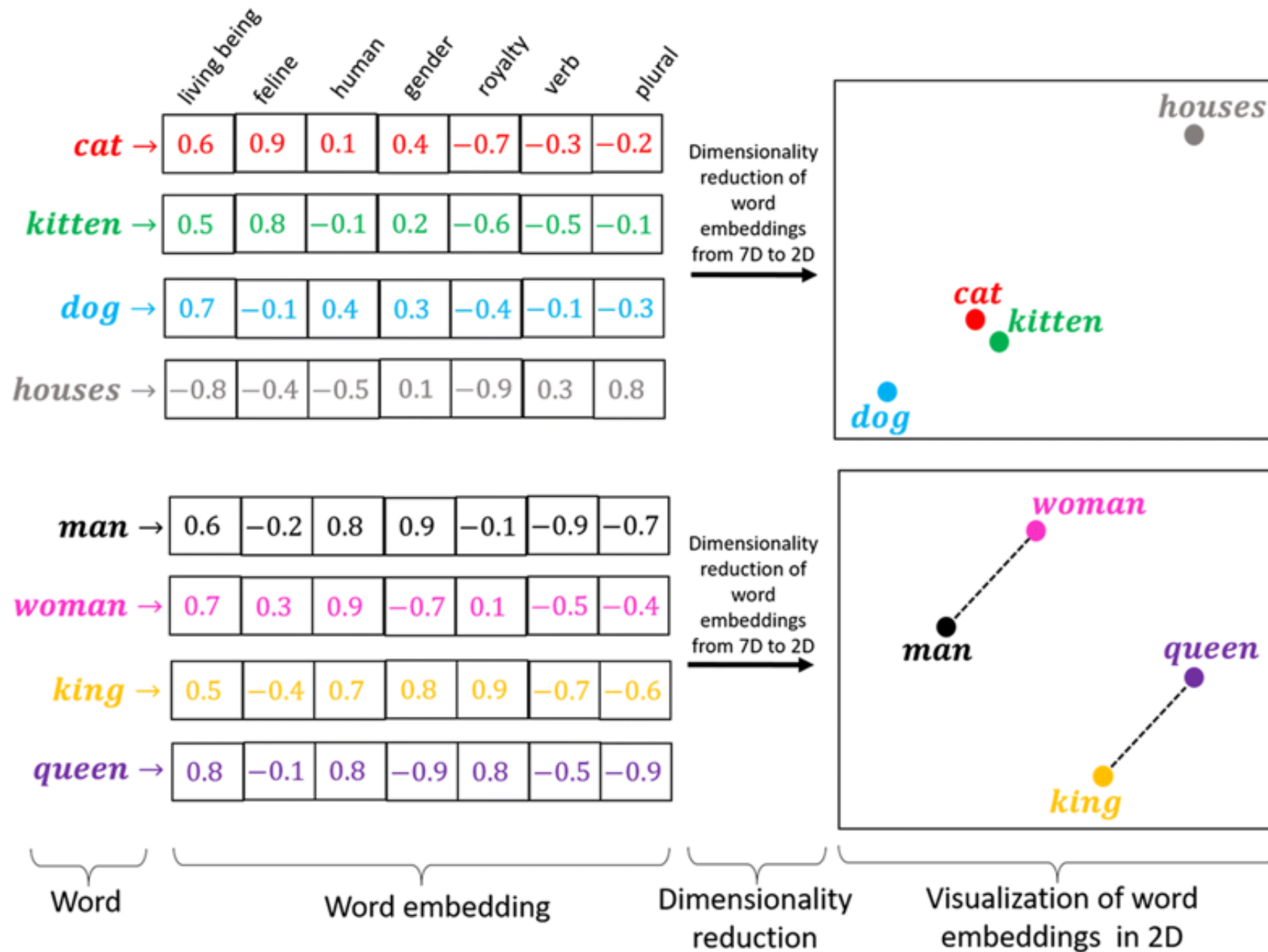
# Distributional approaches

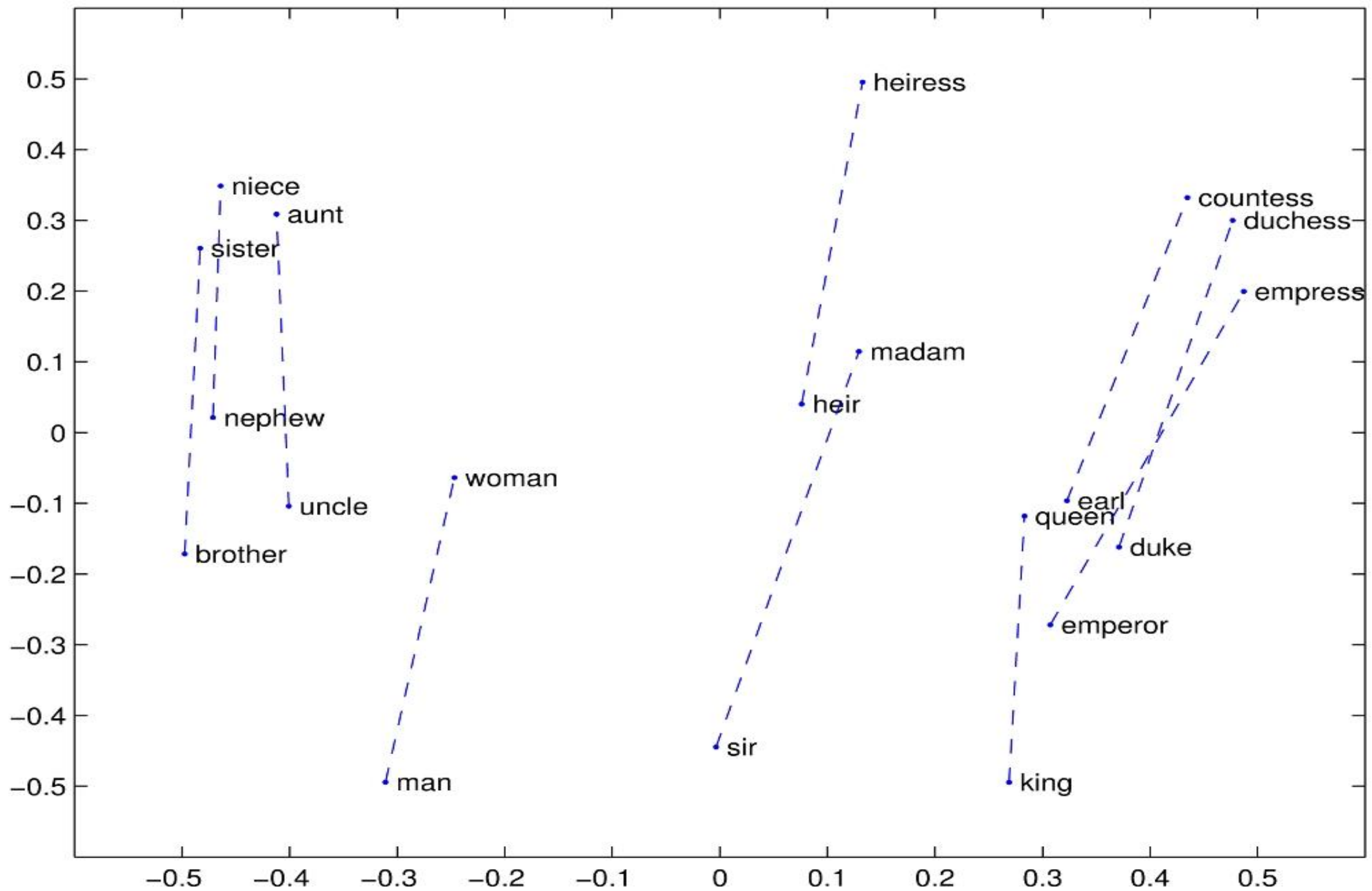- Word meaning is defined by context
  - "You **shall know a word** by the **company it keeps"**
- Context-free word embeddings
  - Word2vec, GloVe, FastText
- Bidirectional approaches
  - ELMo, BERT, etc.

| Word | | | Word embedding | | | | | Dimensionality reduction | Visualization of word embeddings in 2D |
|---|---|---|---|---|---|---|---|---|---|

Word embedding features: living being, feline, human, gender, royalty, verb, plural

| Word | living being | feline | human | gender | royalty | verb | plural |
|---|---|---|---|---|---|---|---|
| cat → | 0.6 | 0.9 | 0.1 | 0.4 | −0.7 | −0.3 | −0.2 |
| kitten → | 0.5 | 0.8 | −0.1 | 0.2 | −0.6 | −0.5 | −0.1 |
| dog → | 0.7 | −0.1 | 0.4 | 0.3 | −0.4 | −0.1 | −0.3 |
| houses → | −0.8 | −0.4 | −0.5 | 0.1 | −0.9 | 0.3 | 0.8 |

Dimensionality reduction of word embeddings from 7D to 2D

| Word | | | | | | | |
|---|---|---|---|---|---|---|---|
| man → | 0.6 | −0.2 | 0.8 | 0.9 | −0.1 | −0.9 | −0.7 |
| woman → | 0.7 | 0.3 | 0.9 | −0.7 | 0.1 | −0.5 | −0.4 |
| king → | 0.5 | −0.4 | 0.7 | 0.8 | 0.9 | −0.7 | −0.6 |
| queen → | 0.8 | −0.1 | 0.8 | −0.9 | 0.8 | −0.5 | −0.9 |

Dimensionality reduction of word embeddings from 7D to 2D

# Similarity can be projected in 2D



not good

to          by

's

dislike          bad          worst

that          now

are

incredibly bad          worse

a          i

you

than          with

is

very good          incredibly good

amazing          fantastic

terrific          wonderful

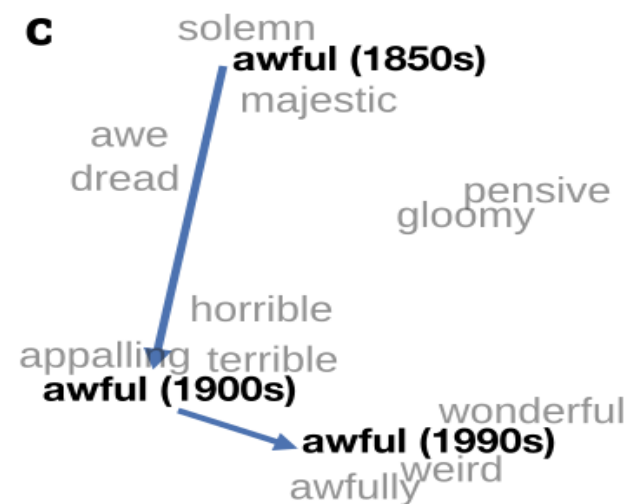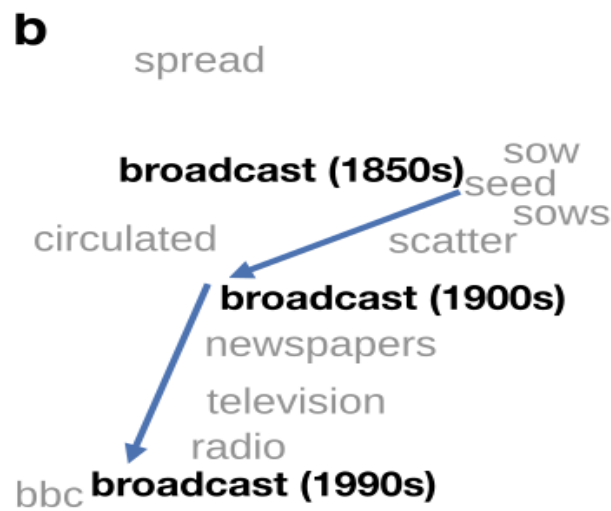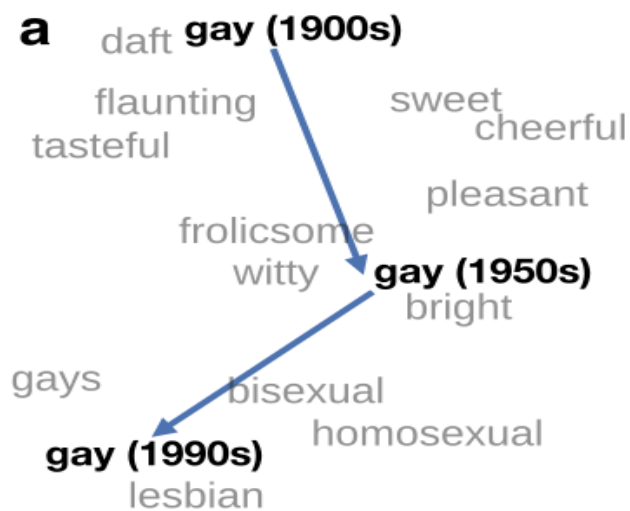nice

good

# Diachronic changes

~30 million books, 1850-1990, Google Books data

# Embeddings reflect cultural biases

- Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. "Man is to computer programmer as woman is to homemaker? debiasing word embeddings." In *Advances in Neural Information Processing Systems*, pp. 4349-4357. 2016.
- Paris : France :: Tokyo : x?
  - x = Japan
- father : doctor :: mother : x?
  - x = nurse
- man : computer programmer :: woman : x?
  - x = homemaker

# Conclusions

- Language is ambiguous!
  - Ambiguity permeates all levels
- NLP has a (relatively) long history
  - Recent years have seen remarkable progress
- NLP applications are all around us
- The tooling, resources and datasets are now much more freely available
  - Universities & textbooks -> Blogs, YouTube, Medium, Coursera, Meetups, Github, etc…

# Thank you

**Tony Russell-Rose, PhD FBCS CITP CEng**
Director, UXLabs | Founder, 2Dsearch
Senior Lecturer, Goldsmiths

- Web:       uxlabs.co.uk, 2dsearch.com
- Email:      tgr@uxlabs.co.uk
- Blog:       http://isquared.wordpress.com
- LinkedIn:  http://uk.linkedin.com/in/tonyrussellrose
- Twitter:    @tonygrr