

# TRI: a tool for the diachronic analysis of large corpora and social media

Incontro utenti ReCaS - Bari, 12 luglio 2019

Pierpaolo Basile  
[pierpaolo.basile@uniba.it](mailto:pierpaolo.basile@uniba.it)



# Hello!

I am Pierpaolo Basile

*Natural Language Processing*

*Distributional Semantics*

*Information Retrieval/Filtering*

You can find me at [pierpaolo.basile@uniba.it](mailto:pierpaolo.basile@uniba.it)

**Words change their  
meaning (*usage*)**

Marty, in 2015  
people will surf on  
the web!!!



# Words change their meaning (*usage*)

Surf!?!?! On  
the  
web!?!?!?



# Motivation

## Detect meaning shift

When was this meaning introduced?

Surf!?!?! On  
the  
web!?!?!?

*surf the Net/Internet to use the Internet*



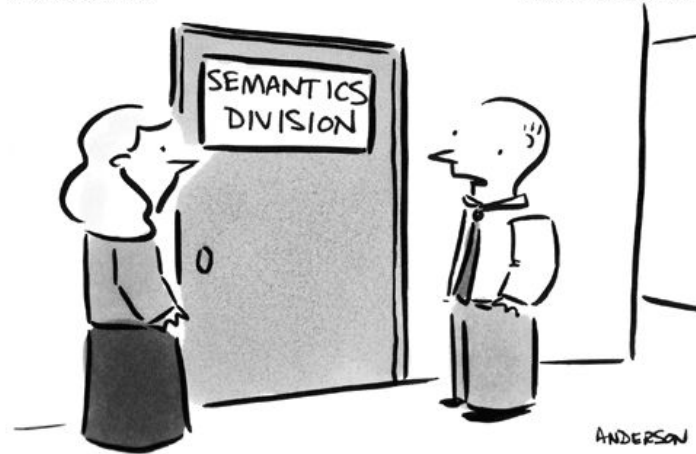
# Diachronic Linguistics

The scientific study of language change over time  
(also called **Historical Linguistics**)

# How to represent semantics?

© MARK ANDERSON

WWW.ANDERSTOONS.COM



ANDERSON

"We're really more of a department."

# Distributional Semantic Models

- Analysis of word-usage statistics over huge corpora
- Geometric space of concepts
- Similar words are represented close in the space

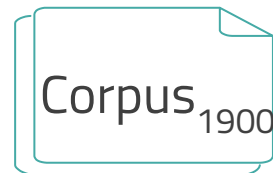
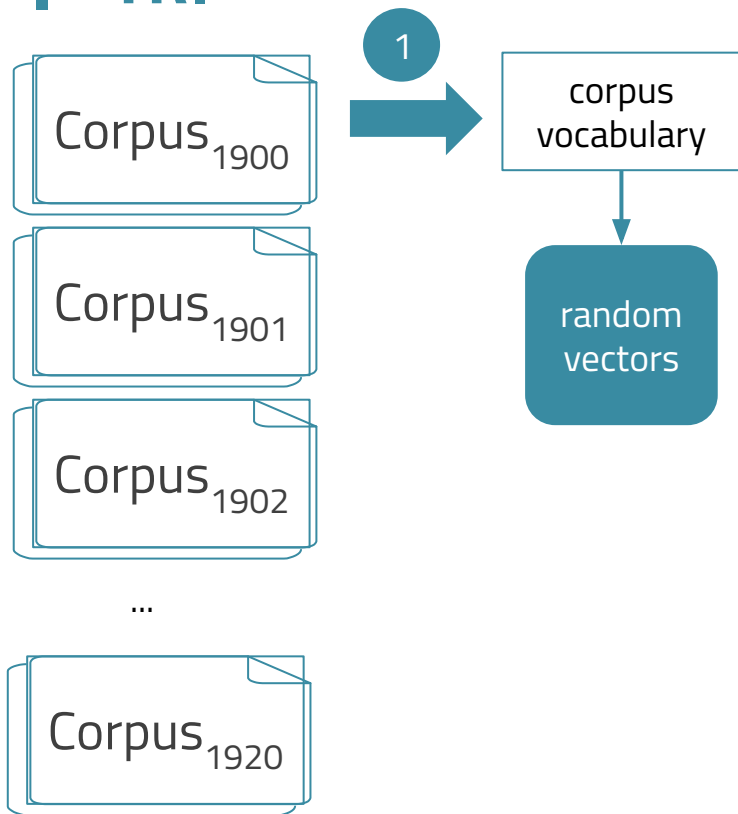
memory floppy\_disk  
ram chip disk hard\_disk  
software printer  
computer  
os workstation  
pc device  
operating\_system  
linux mouse  
tux mice  
penguin rabbit rat  
dog animal  
cat monkey insect



“ A **WordSpace** is a snapshot of a specific corpus, it does not take into account temporal information

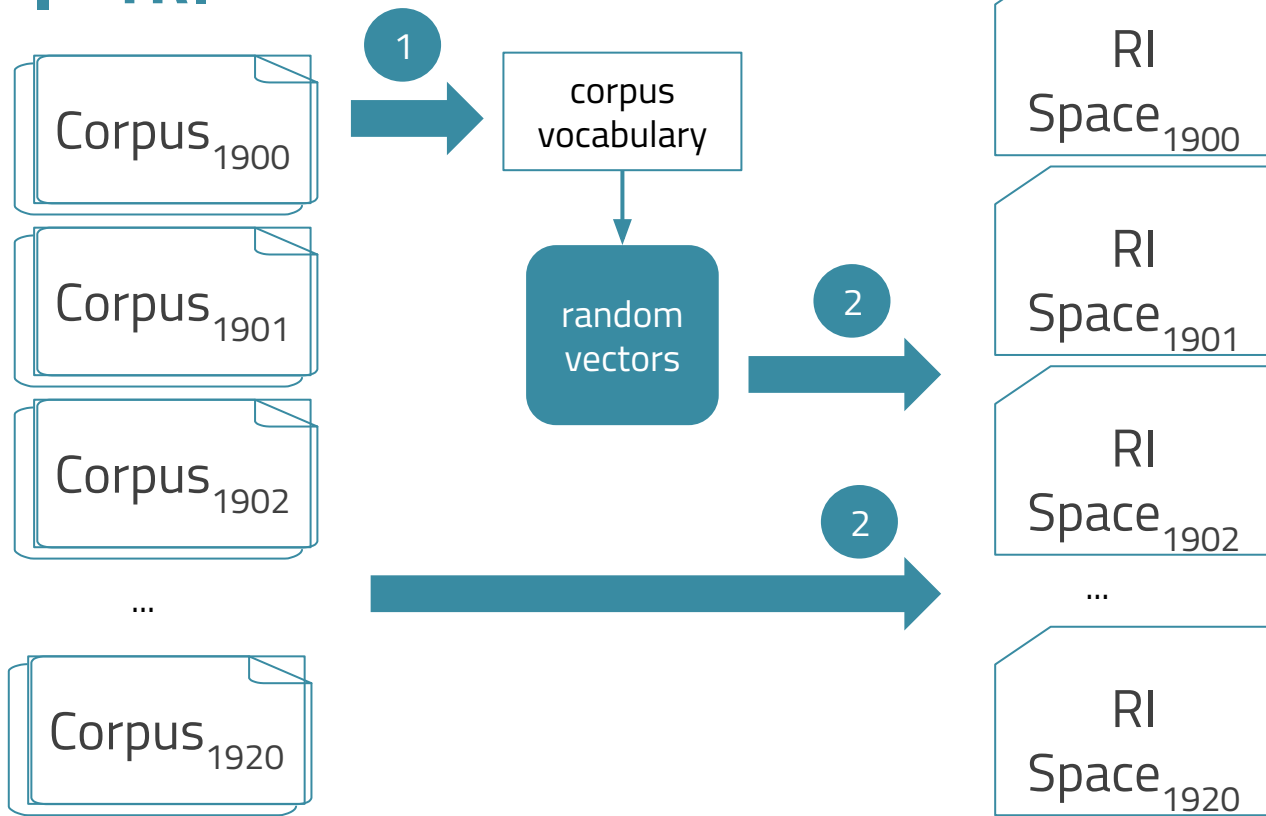
# Temporal Random Indexing

## TRI



# Temporal Random Indexing

## TRI



**Semantic vector** for a term is the sum of the context vectors co-occurring with the term in the same time period

# Temporal Random Indexing

## TRI

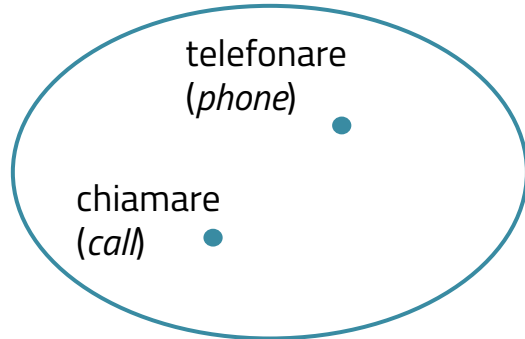
- Corpus with temporal information
  - split the corpus in several time periods
- Build a WordSpace for each time period using TRI
- Words in different WordSpaces are **comparable!**

# Similarity between words can change over time

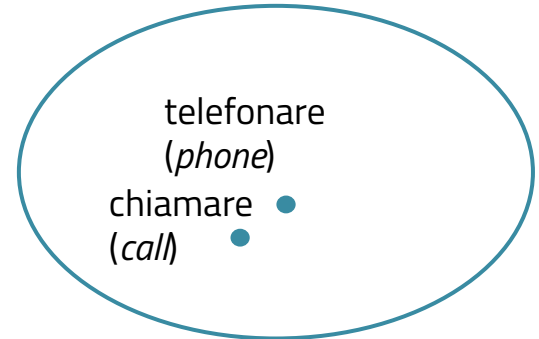
WordSpace 1870



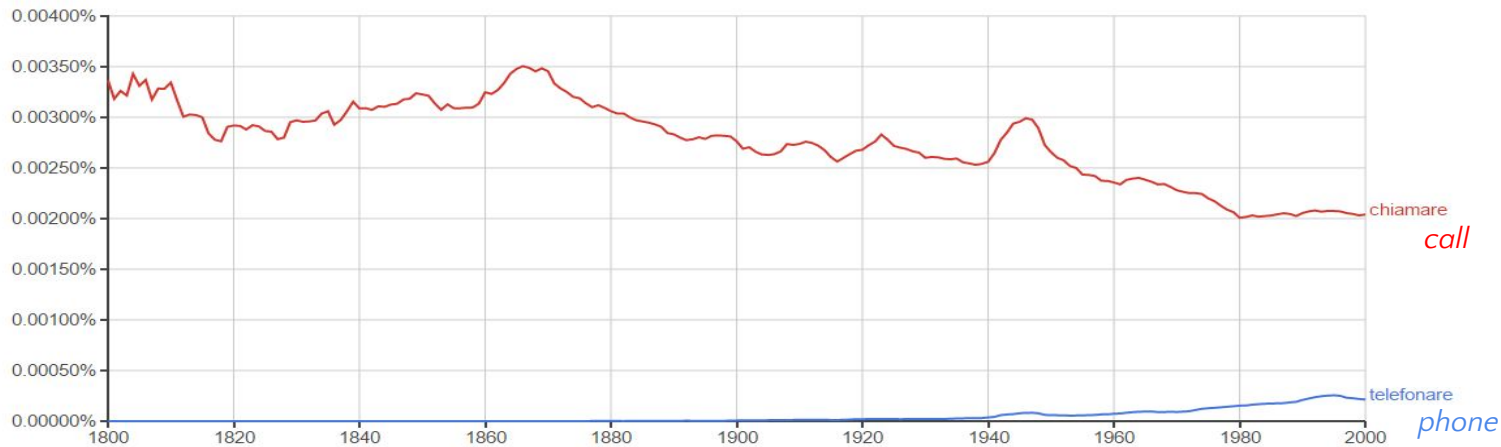
WordSpace 1920



WordSpace 1930



# Google N-gram



# TRI



# Methodology



**TRI**

Run TRI on a corpus split in time periods

**Time Series**

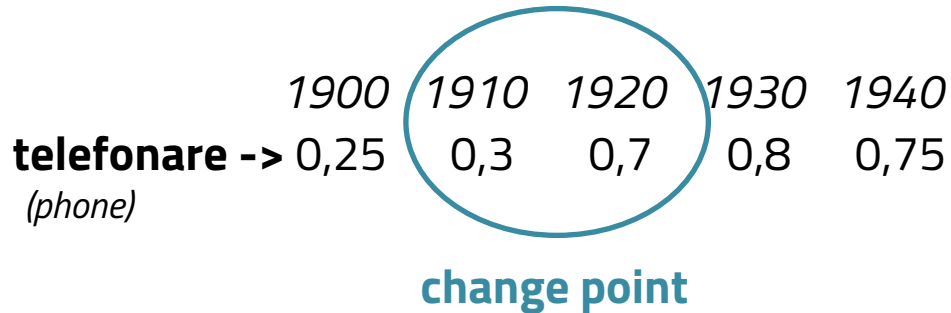
Provide a time series for each word

**Change Point Detection**

Detect significant changes in the time series

# Change point detection

- Track the word meaning change over time
- Build a time series by taking into account the semantic shift of each word
- Find significant change: **Mean shift model**





# Social media

- Build TRI on **Twitter**
- About **500M tweets** (feb. 2012 – sep. 2015)
- Time interval = **1 month**
- Change point detection on the 1,000 most frequent hashtags

# Social media

## #bologna (august 2014)



### #bologna (august 2014)

vittime

bologna

memoria

strage

chiedere

festa

corsa

manifestazione

# Social media

## #euro (giugno 2015)

LA CRISI ELLENICA

### Il parlamento greco approva il referendum. Tsipras chiede di votare «no»

—con un'artcle gallery di **Vittorio Da Rold** e un post di Econopoly | 28 giugno 2015

#euro (june 2015)

#europa

#grexit

ora

**alternativa**

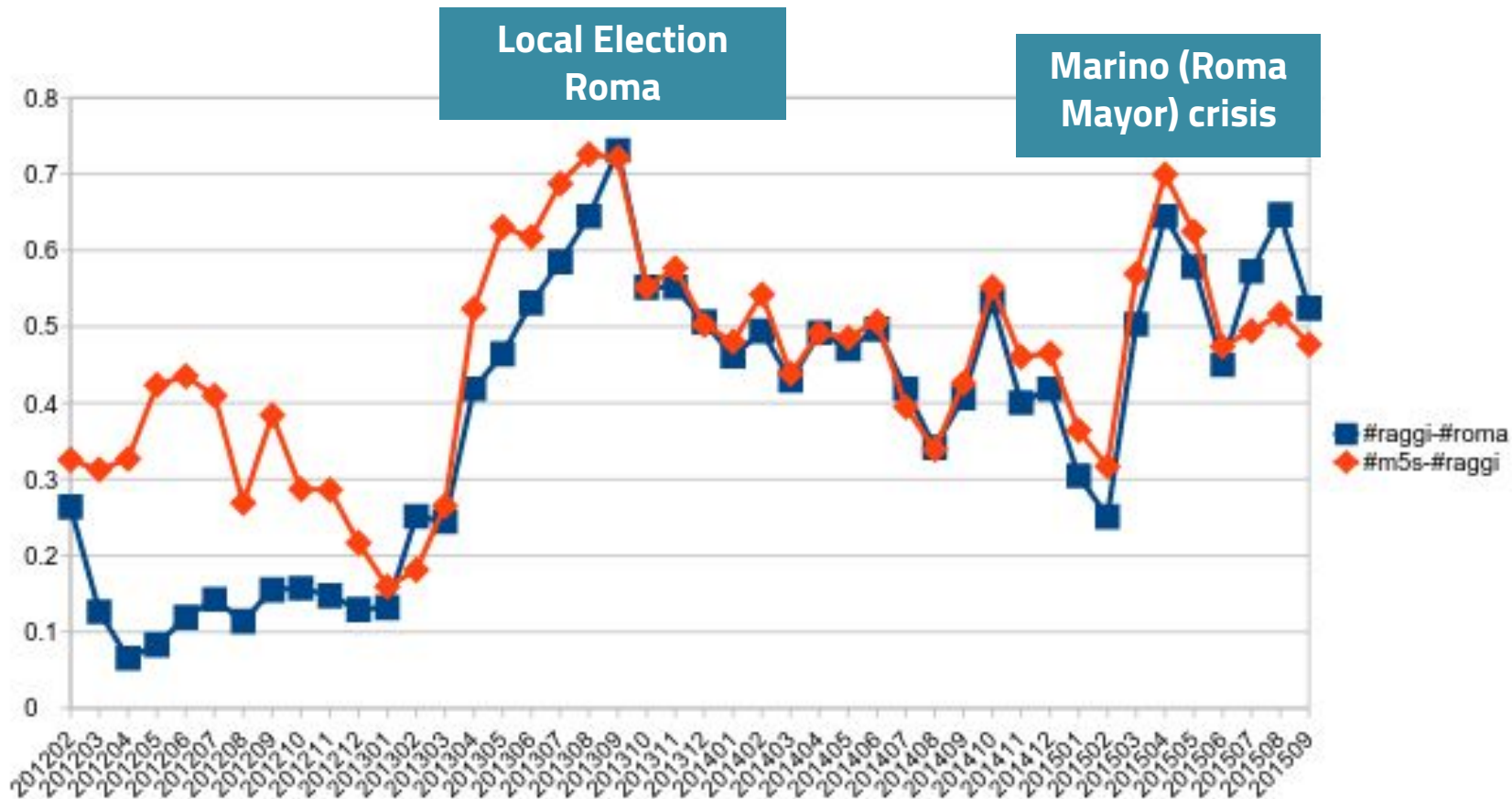
acqua

invasione

**uscita**

economia

# Social media



# Ongoing work

# Build a gold standard for the evaluation

## Dizionario di Italiano

*il Sabatini Coletti* Dizionario della Lingua Italiana

[Codice da incorporare »](#)

CERCA

Dizionario di Italiano

girocollo  
giroconto  
giromanica  
girondino  
girone  
gironzolare  
giropilota  
niroscnico

**girotondo** [gi-ro-tón-do] s.m. inv.

- 1 Gioco infantile consistente nel formare un cerchio tenendosi per mano e nel girare cantando una filastrocca
- 2 Manifestazione politica di protesta non organizzata da partiti

• a. 1869 (1); a. 2001 (2)

change point

# Evaluation

- Build word vectors using different approaches exploiting the Italian Google n-grams corpus
  - TRI, word embeddings alignment
- Evaluation using the gold standard
  - time period: 1900-2012

# Thanks!!

## Any questions?

You can find me at  
[pierpaolo.basile@uniba.it](mailto:pierpaolo.basile@uniba.it)