

# Algorithms for Knowledge and Information Extraction in Text with Wikipedia



**Marco Ponza**

University of Pisa

Supervisor

Prof. Paolo Ferragina



# Menu



## 0. Introduction

### Knowledge and Information Extraction

1. Algorithms for Entity Relatedness
2. Algorithms for Entity and Fact Salience

### Applications

3. Algorithms for Expert Finding
4. Future Research Directions

0

# Introduction

# Introduction

- ▷ Enhancing the humankind progress with new **intelligent** technologies



Tools that can afford general- or specific-domain **tasks** with **performance close** or **better** than humans

- ▷ Machines need of access, read and **understand** information stored in **data archives**



The **dominant** form on which information is produced every day by humans is still **Natural Language**

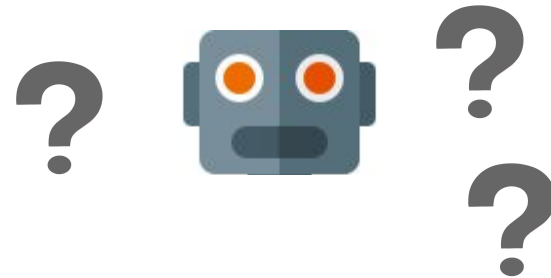
# Introduction

## Text Understanding

Easy for humans



Hard for machines



# Introduction

## Text Understanding



1

Map ambiguous words into the **real-world entities** they refer to as well as **contextualize** them together with **related** entities

# Introduction

## Text Understanding

*Leonardo is the scientist who painted Mona Lisa*



(“Leonardo”, “is”, “scientist”)

(“Leonardo”, “painted”, “Mona Lisa”)

2

Structure multiple **facts** (propositions)  
contained in the sentence

Triples of (**subject**, **relation**, **object**)

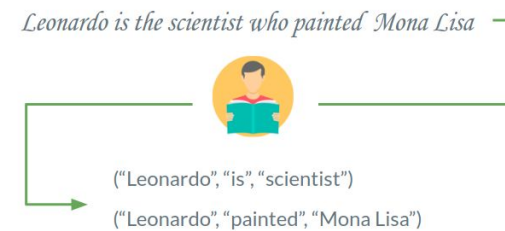
# Introduction

## Text Understanding



1

Map words into the **concepts** they refer to



2

Structure multiple **facts** (propositions) contained in the sentence  
Triplets of (subject, relation, object)

How can we do that?

Humans can interpret words in a *larger context* hinging onto their **background** and **linguistic knowledge** (Gabrilovich, SIGIR'16)

Detect (1) unambiguous **entities** (2) **facts**, (3) quantifying how much they are **related**, and (4) **efficiently retrieve** related entities



# Introduction

## Text Understanding



- ▷ Literature currently offers a number of solutions based on **BoW** (Harris, Word'54): *A text is a vector of ambiguous keywords*

### Limitations

- Curse of Dimensionality
- Synonymy and Polysemy problems of keywords
- No understanding of real-world entities
- Structure of the sentence is lost (no facts)

# Introduction

## Text Understanding



- ▷ Literature currently offers a number of solutions based on **BoW** (Harris, Word'54) : *A text is a vector of ambiguous keywords*
  - Limitations** {
    - Curse of Dimensionality
    - Synonymy and Polysemy problems of keywords
    - No understanding of real-world entities
    - Structure of the sentence is lost (no facts)
- ▷ **LDA/LSI** (Huffman, NIPS'10) and **Word Embeddings** (Mikolov, NIPS'13) overcome some limitations
  - A text is mapped into a **latent space** (vector of floating-points)*

# Introduction

## Text Understanding



- ▷ Literature currently offers a number of solutions based on **BoW** (Harris, Word'54) : *A text is a vector of ambiguous keywords*

### Limitations

- Curse of Dimensionality ✓
- Synonymy and Polysemy problems of keywords ✓
- No understanding of real-world entities ✗ ✓
- Structure of the sentence is lost (no facts) ✗

- ▷ **LDA/LSI** (Huffman, NIPS'10) and **Word Embeddings** (Mikolov, NIPS'13) overcome some limitations

*A text is mapped into a latent space (vector of floating-points)*

# Introduction

## Text Understanding



- ▶ Literature currently offers a number of solutions based on **BoW** (Harris, Word'54) : *A text is a vector of ambiguous keywords*

Need for a more efficient and effective semantic paradigm

- ▶ **LDA/LSI** (Huffman, NIPS'10) and **Word Embeddings** (Mikolov, NIPS'13) overcome some limitations

A text is mapped into a *latent space* (vector of floating-points)

# Introduction

- ▷ Need for a more efficient and effective semantic paradigm

Exploiting two  
different resources



- World Knowledge

**WIKIPEDIA**  
The Free Encyclopedia

- Linguistic Knowledge

Language Grammar

...thanks to recent advancements in the field of  
Natural Language Processing:

- ▷ **Entity Linking** (Bunescu, EACL'06), (Scaiella, IEEE'12), (Piccinno, SIGIR'14)

# Introduction

- ▷ Need for a more efficient and effective semantic paradigm

Exploiting two  
different resources



- World Knowledge

WIKIPEDIA  
The Free Encyclopedia

- Linguistic Knowledge

Language Grammar

...thanks to recent advancements in the field of  
Natural Language Processing:

- ▷ **Entity Linking** (Bunescu, EACL'06), (Scaiella, IEEE'12), (Piccinno, SIGIR'14)

*Leonardo painted Mona Lisa*

# Introduction

- ▷ Need for a more efficient and effective semantic paradigm

Exploiting two  
different resources

- World Knowledge

WIKIPEDIA  
The Free Encyclopedia

- Linguistic Knowledge

Language Grammar

...thanks to recent advancements in the field of  
Natural Language Processing:

- ▷ Entity Linking (Bunescu, EACL'06), (Scaiella, IEEE'12), (Piccinno, SIGIR'14)

*Leonardo painted Mona Lisa*



Leonardo da Vinci



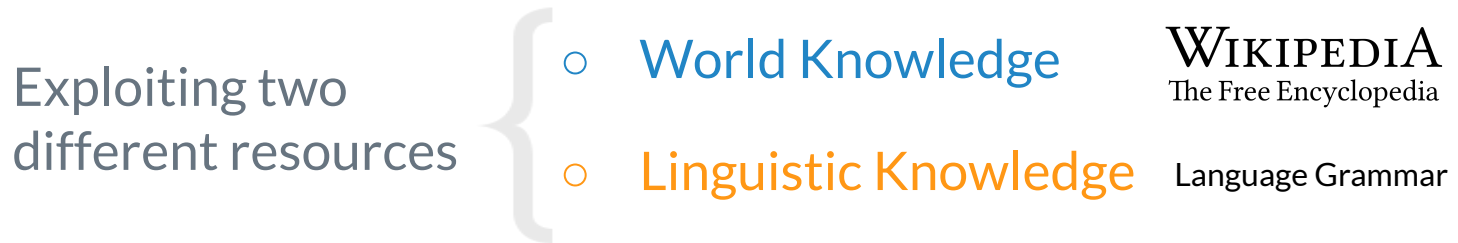
Mona Lisa (painting)

Entities



# Introduction

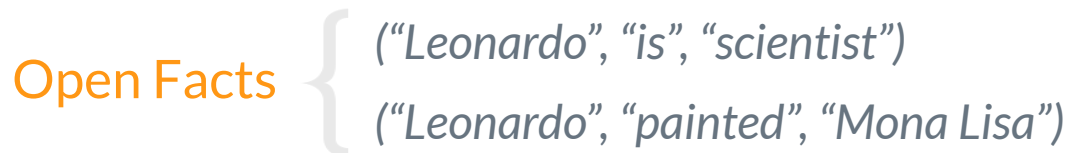
▷ **Need for a more efficient and effective semantic paradigm**



...thanks to recent advancements in the field of Natural Language Processing:

- ▷ **Entity Linking** (Bunescu, EACL'06), (Scaiella, IEEE'12), (Piccinno, SIGIR'14)
- ▷ **Open Information Extraction** (Banko, IJCAI'07),(Del Corro, WWW'13), (Gashteovski, EMNLP'17)

*Leonardo, the scientist, painted Mona Lisa*



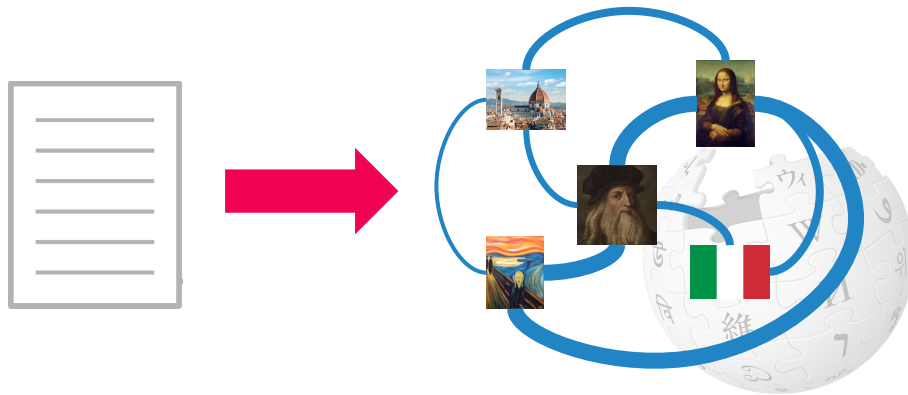


# Introduction

- ▷ Need for a more efficient and effective semantic paradigm

...how?

Applying **Graph Theory** to entity linking and open information extraction



Model a text as a ***Small Wikipedia Graph!***

- Curse of Dimensionality ✓ →
- Synonymy and Polysemy problems of keywords ✓ →
- Understanding of real-world entities ✓ →
- Structured facts ✓ →

The graph is small

Wikipedia entities are unique and they represent a real-world concept

OpenIE preserves subject-relation-object structure

# Contributions

1 Entity Relatedness

How much two Wikipedia entities are *related*?

2 Entity Saliency

*Summarize* document's subject matter with its *Salient Wikipedia Entities*

3 Fact Saliency

*Summarize* document's subject matter with its *Salient Open Facts*

4 Expert Finding

Who are the *experts* of a given topic?



# Contributions

1

## Entity Relatedness

A Two-Stage Framework for Computing  
Entity Relatedness in Wikipedia  
Marco Ponza, Paolo Ferragina and Soumen Chakrabarti



2

## Entity Salience

Document Aboutness via Sophisticated  
Syntactic and Semantic Features  
Marco Ponza, Paolo Ferragina and Francesco Piccinno



SWAT: A System for Detecting  
Salient Wikipedia Entities in Texts  
Marco Ponza, Paolo Ferragina and Francesco Piccinno



3

## Fact Salience

Facts That Matter  
Marco Ponza, Luciano Del Corro and Gerhard Weikum



4

## Expert Finding

WISER: A Semantic Approach for Expert Finding  
in Academia based on Entity Linking  
Paolo Cifariello, Paolo Ferragina and Marco Ponza



Information Systems 2019

# 1

# Algorithms for Entity Relatedness

A Two-Stage Framework for Computing  
Entity Relatedness in Wikipedia

Marco Ponza, Paolo Ferragina and Soumen Chakrabarti



# Entity Relatedness

## Motivation

Proliferation of the usage of **Knowledge Graphs**



Freebase



Consumers

- ▷ **Retrieval of Information** (Blanco, WSDM '15), (Cornolti, WWW '16)
- ▷ **Entity Linking** (Mihalcea, CIKM '07), (Meij, WSDM '12), (Ganea, WWW '16)
- ▷ **Document Clustering, Classification and Similarity**  
(Scaiella, WSDM '12), (Vitale, ECIR '12), (Ni, WSDM '16)



Need for computing *entity relatedness*

Compute how much two **entities** are **related**

*Relatedness* : **Entities** × **Entities** → **Real**



# The Wikipedia Knowledge Graph

- ▶ Our Knowledge Graph (KG): **WIKIPEDIA**  
The Free Encyclopedia



# The Wikipedia Knowledge Graph

- ▶ Our Knowledge Graph (KG): **WIKIPEDIA**  
The Free Encyclopedia
  - Entity?





**WIKIPEDIA**  
The Free Encyclopedia

[Main page](#)  
[Contents](#)  
[Featured content](#)  
[Current events](#)  
[Random article](#)  
[Donate to Wikipedia](#)  
[Wikipedia store](#)

**Interaction**  
[Help](#)  
[About Wikipedia](#)  
[Community portal](#)  
[Recent changes](#)  
[Contact page](#)

**Tools**  
[What links here](#)  
[Related changes](#)  
[Upload file](#)  
[Special pages](#)  
[Permanent link](#)  
[Page information](#)  
[Wikidata item](#)  
[Cite this page](#)

Article [Talk](#)

Read [View source](#) [View history](#)

# Leonardo da Vinci



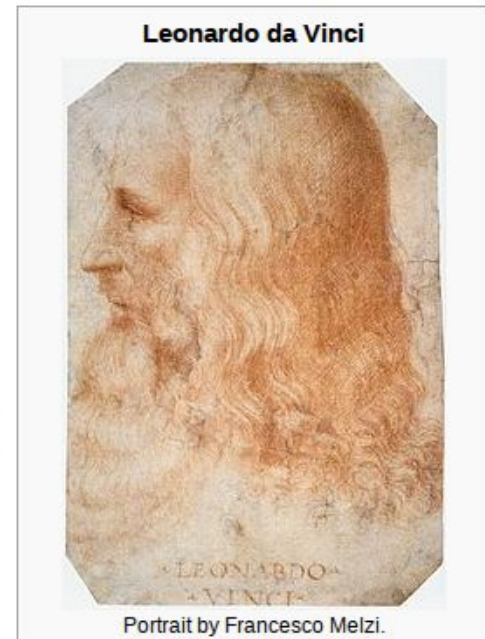
From Wikipedia, the free encyclopedia

*"Da Vinci" redirects here. For other uses, see [Da Vinci \(disambiguation\)](#).*

*This is a [Renaissance Florentine](#) name. The name *da Vinci* is an indicator of birthplace, not a family name; this person is properly referred to by the given name *Leonardo*.*

**Leonardo di ser Piero da Vinci** (Italian: [leoˈnardo di ˌsɛr ˈpjɛːro da (v)ˈvintʃi] ⓘ<sup>ⓘ</sup> listen<sup>ⓘ</sup>), more commonly **Leonardo da Vinci** or simply **Leonardo** (15 April 1452 – 2 May 1519), was an Italian polymath whose areas of interest included invention, painting, sculpting, architecture, science, music, mathematics, engineering, literature, anatomy, geology, astronomy, botany, writing, history, and cartography. He has been variously called the father of palaeontology, ichnology, and architecture, and is widely considered one of the greatest painters of all time. Sometimes credited with the inventions of the parachute, helicopter and tank,<sup>[1][2][3]</sup> he epitomised the Renaissance humanist ideal.

Many historians and scholars regard Leonardo as the prime exemplar of the "[Universal Genius](#)" or "Renaissance Man", an individual of "unquenchable curiosity" and "feverishly inventive imagination".<sup>[4]</sup> According to art historian [Helen Gardner](#), the scope and depth of his interests were without precedent in recorded history, and "his mind and personality seem to us superhuman, while the man himself mysterious and remote".<sup>[4]</sup> Marco Rosci notes that while there is much speculation regarding his life and personality, his view of the world was logical rather than mysterious, and that the empirical methods he employed were unorthodox for his time.<sup>[5]</sup>



▷ **Entity** = Wikipedia Page = Node of our KG





# Terminology

▷ Our Knowledge Graph (KG):

- **Entity** = Wikipedia Page (a node of KG)
- **Label** = Textual Description of the Wikipedia Page
- **Edges?**

# WIKIPEDIA

The Free Encyclopedia





WIKIPEDIA  
The Free Encyclopedia

- Main page
- Contents
- Featured content
- Current events
- Random article
- Donate to Wikipedia
- Wikipedia store

- Interaction
  - Help
  - About Wikipedia
  - Community portal
  - Recent changes
  - Contact page

- Tools
  - What links here

**Science**

From Wikipedia, the free encyclopedia

*This article is about the general term. For other uses, see Science (disambiguation).*

**Science** is a systematic enterprise that builds and organizes knowledge in the form of testable explanations and predictions about the universe.

Contemporary science is typically subdivided into the natural sciences, which study the material universe; the social sciences which study people and societies; and the formal sciences, such as mathematics. The formal sciences are often excluded as they do not depend on empirical observations. Disciplines which use science like engineering and medicine may also be considered to be applied sciences.

During the Middle Ages in the Middle East, foundations for the scientific method were laid by Ibn al-Haytham in his *Book of Optics*. From classical antiquity through the 19th century, science as a type of knowledge was more closely linked to philosophy than it is now and, in fact, in the Western world, the term "natural philosophy" encompassed fields of study that are today associated with science, such as astronomy, medicine, and physics. While the classification of the material world of the ancient Indians and Greeks into air, earth, fire and water was more philosophical, medieval Middle Eastern scientists used practical, experimental observation to classify materials.

In the 17th and 18th centuries, scientists increasingly sought to formulate knowledge in terms of laws of nature. Over the course of the 19th century, the word "science" became increasingly associated with the scientific method itself, as a disciplined way to study the natural world. It was in the 19th century that scientific disciplines such as biology, chemistry, and physics reached their modern shapes. The same time period also included the origin of the terms "scientist" and "scientific community," the founding of scientific institutions, and increasing significance of the interactions with society and other aspects of culture.

- Wikimedia Commons
- Wikiquote
- Wikisource

- Languages
  - Адыгэбзэ

Article [Talk](#)

Read [View source](#) [View history](#)

Search

# Leonardo da Vinci

From Wikipedia, the free encyclopedia

*"Da Vinci" redirects here. For other uses, see Da Vinci (disambiguation).*

*This is a Renaissance Florentine name. The name da Vinci is an indicator of birthplace, not a family name; this person is properly referred to by the given name Leonardo.*

**Leonardo di ser Piero da Vinci** (Italian: [leoˈnardo di ˌsɛr ˈpjɛːro da (v)ˈvintʃi] (listen)), more commonly **Leonardo da Vinci** or simply **Leonardo** (15 April 1452 – 2 May 1519), was an Italian polymath whose areas of interest included invention, painting, sculpting, architecture, science, music, mathematics, engineering, literature, anatomy, geology, astronomy, botany, writing, history, and cartography. He has been variously called the father of palaeontology, ichnology, and architecture, and is widely considered one of the greatest painters of all time. Sometimes credited with the inventions of the parachute, helicopter and tank,<sup>[1][2][3]</sup> he epitomised the Renaissance humanist ideal.

Many historians and scholars regard him as a "Genius" or "Renaissance Man".

**Invention**

From Wikipedia, the free encyclopedia

*"Inventor" and "Invented" redirect here. For other uses, see Invention (disambiguation).*

*For more details on inventions throughout history, see Timeline of historic inventions.*

*For the CAD design software, see Autodesk Inventor.*

An **invention** is a unique or novel device, method, composition or process. The invention process is engineering and product development process. It may be an improvement upon a machine, process or system, or a completely new object or a result. An invention that achieves a completely unique function or result may be novel and **not obvious to others skilled in the same field**. An inventor may be taking a big step by improving upon a preexisting idea or invention, or by coming up with something wholly new. Some inventions can be patented. A patent legally protects the intellectual property rights of an inventor. The rules and requirements for patenting an invention process of obtaining a patent is often expensive.

Another meaning of invention is **cultural invention**, which is an innovative set of useful social practices or ideas passed on to others.<sup>[1]</sup> The Institute for Social Inventions collected many such ideas in magazines and books. An important component of artistic and design creativity. Inventions often extend the boundaries of human capability.



**Astronomy**

From Wikipedia, the free encyclopedia

*This article is about the scientific study of celestial objects. For other uses, see Astronomy (disambiguation).*

**Astronomy**, a natural science, is the study of celestial objects (such as stars, galaxies, planets, moons, and nebulae) and processes (such as supernovae explosions, gamma ray bursts, and cosmic microwave background radiation), and evolution of such objects and processes, and more generally all phenomena in the universe. A related but distinct subject, physical cosmology, is concerned with studying the universe as a whole. Astronomy is the oldest of the natural sciences. The early civilizations in recorded history, such as the Egyptians, Nubians, Iranians, Chinese, and Maya performed methodical observations of the night sky. The included disciplines as diverse as astrometry, celestial navigation, observational astronomy and professional astronomy is nowadays often considered to be synonymous with astrophysics.<sup>[2]</sup>

During the 20th century, the field of professional astronomy split into observational and theoretical astronomy. Theoretical astronomy is focused on acquiring data from observations of astronomical objects, which is then analyzed using mathematical models. Theoretical astronomy is oriented toward the development of computer or analytical models to explain observations being used to confirm theoretical results.

Astronomy is one of the few sciences where amateurs can still play an active role, especially in the discovery of transient phenomena. Amateur astronomers have made and contributed to many important astronomical discoveries, including new comets.

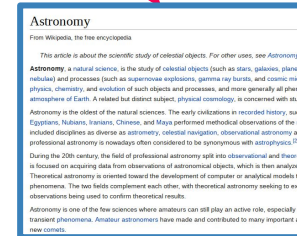
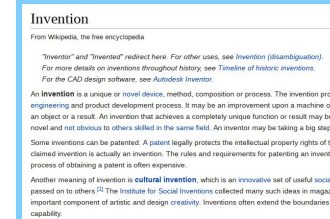
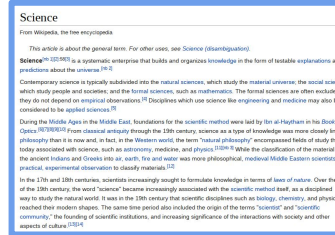
Leonardo was, and is, renowned primarily as a painter. Among his works, the *Mona Lisa* is his most famous and most parodied portrait<sup>[6]</sup> and *The Last Supper* the most reproduced religious painting of all time, their fame approached only by Michelangelo's *The Creation of Adam*. Leonardo's drawing of the *Vitruvian Man* is also regarded as a cultural icon,<sup>[7]</sup> being

*The Vitruvian Man*

# The Wikipedia Knowledge Graph

## ▷ Our Knowledge Graph (KG): **WIKIPEDIA** The Free Encyclopedia

- **Entity** = Wikipedia Page (a node of KG)
- **Label** = Textual Description of the Wikipedia Page
- **Edge** = Wikipedia Hyperlinks

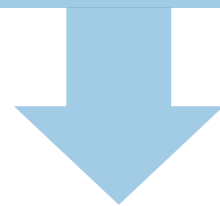


# Known Relatedness Methods

A large number of methods proposed in literature...

- **Personalized Web Search** (Haveliwala, WWW '02)
- **Link Prediction** (Liben-Nowell, JAIST '07)
- **Word and Document Similarity** (Gabrilovich, IJCAI '07)
- **Document Annotation** (Piccinno, SIGIR '14)
- **Machine Translation** (Rothe, ACL '14)
- **Document Classification** (Perozzi, KDD '14), (Tan, WWW '15)

...that have been applied or are similar to our problem



We have experimented them  
on the Entity Relatedness task

# Introduction

## Our Contributions

1

### New dataset WiRe

- Human-assigned scores
- 503 Wikipedia entity pairs
- Sampled from New York Times (Dunietz, EACL '14)

2

### Thorough and systematic study of *all known relatedness measures*

- WiRe (our introduced dataset)
- WikiSim (Milne, AAAI '08)

3

### Proposal of a **Two-Stage Framework**

- Space-efficient
- Computationally lightweight
- More accurate than previous proposals

4

### Extrinsic evaluation of our proposal

- Domain of **Entity Linking**
- Increase of accuracy and robustness of **TAG-ME** (Scaiella, CIKM '10)

5

Publicly available **WiRe dataset** and the **code of all algorithms!**



# Our Two-Stage Framework

- ▷ Built on the top of existing relatedness algorithms
- ▷ **Improves** current approaches
  - More **accurate** relatedness scores
  - **Fast** at query time
- ▷ The two stages of our framework:

1

A **small** and **weighted subgraph** is dynamically grown around the two *query entities*

2

Computing the **relatedness** between the two *query entities* according with the generated subgraph

## ▷ Motivations

- Wikipedia **edges** are **noisy** (introduced for **citation, explanation, ...**)
- Subgraph **nodes** are **strongly related** to the query entities (they are good bridges)
- Subgraph **edges** are **less noisy** (confined to few meaningful bridge nodes)

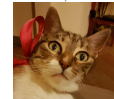
# Our Two-Stage Framework

1

A **small** and **weighted subgraph** is dynamically grown around the two *query entities*



Tiger



Cat





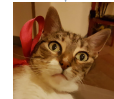
# Our Two-Stage Framework

1

A **small** and **weighted subgraph** is dynamically grown around the two *query entities*



Tiger



Cat

How can we populate the subgraph?

# Our Two-Stage Framework

1

A **small** and **weighted subgraph** is dynamically grown around the two *query entities*



Tiger



Siberian\_tiger



Leopard



Jaguar



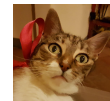
European\_cat



Cat\_anatomy



Felidae



Cat

*Populating the subgraph.* Choosing the **top-k** nodes most related to the query entities

# Our Two-Stage Framework

1

A **small** and **weighted subgraph** is dynamically grown around the two *query entities*

How?



Siberian\_tiger



European\_cat

Various algorithms:

- ESA (Gabrilovich, IJCAI '07)
- Milne-Witten (Milne, AAAI '08)
- DeepWalk (Perozzi, KDD '14)
- Entity2Vec (Ni, WSDM '16)



Leopard



Cat anatomy



cat



Jaguar



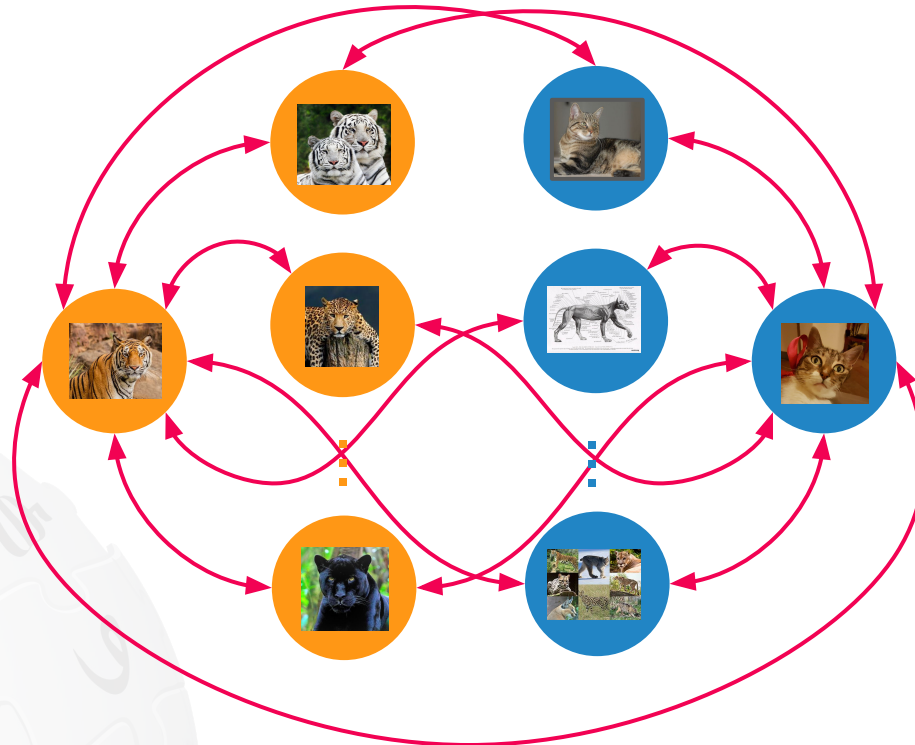
Felidae

*Populating the subgraph.* Choosing the **top-k** nodes **most related** to the query entities

# Our Two-Stage Framework

1

A **small** and **weighted subgraph** is dynamically grown around the two *query entities*



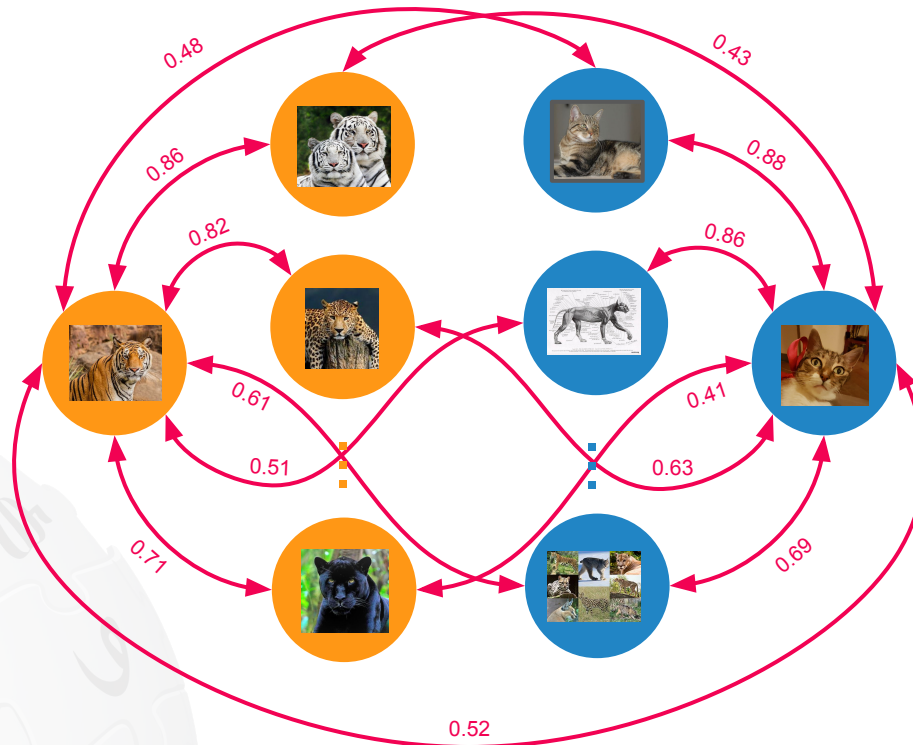
Creating the edges. Each query entity is linked to

- the *other* query entity
- its top-k related entities
- the other top-k related entities

# Our Two-Stage Framework

1

A **small** and **weighted subgraph** is dynamically grown around the two *query entities*



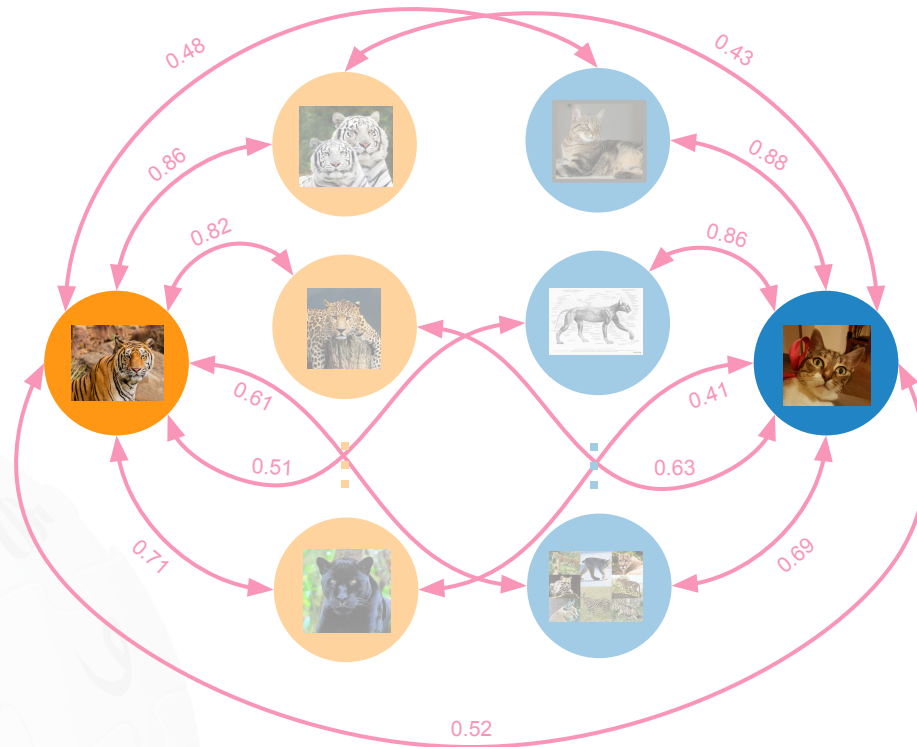
*Weighting the edges.* How?

- Milne-Witten (Milne, AAAI '08)
- DeepWalk (Perozzi, KDD '14)
- Entity2Vec (Ni, WSDM '16)

# Our Two-Stage Framework

2

Computing the **relatedness** between the two *query entities* according with the generated subgraph



Computing Relatedness

**CoSimRank** (Rothe, ACL '14)

$$\text{relatedness} \left( \text{tiger}, \text{cat} \right) = 0.65$$

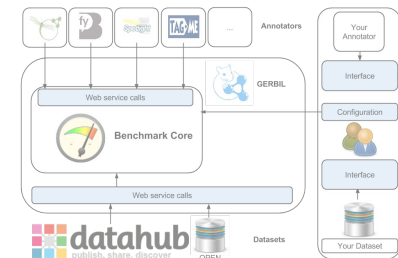
# Experiments

## ▷ Intrinsic evaluation on pairs of Wikipedia Entities

Dataset	WikiSim (Milne, AAAI '08)	WiRe
Size	268	503
Pair Type	Common Nouns	Named Entities
Ground-Truth	Crowdsourcing	Human Experts

## ▷ Extrinsic evaluation

- Domain of Entity Linking
- On four different datasets (Usbeck, WWW '15)



# Experiments

## Intrinsic Evaluation

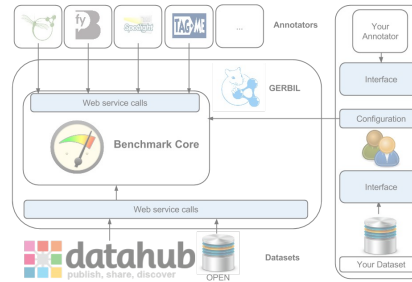
Method	WikiSim			WiRe			AVG
	Pearson	Spearman	Harmonic	Pearson	Spearman	Harmonic	
ESA	0.61	0.72	0.67	0.60	0.63	0.62	0.645
Milne-Witten	0.62	0.65	0.63	0.77	0.69	0.72	0.675
DeepWalk	0.71	0.70	0.71	0.74	0.68	0.71	0.710
Entity2Vec	0.68	0.70	0.69	0.74	0.70	0.72	0.705
<b>Two-Stage Framework</b>	<b>0.74</b>	<b>0.75</b>	<b>0.74</b>	<b>0.83</b>	<b>0.75</b>	<b>0.79</b>	<b>0.765</b>

- ▷ Pearson measures predicted-vs-correct scores
- ▷ Spearman focuses on the ranking order among entity pairs
- ▷ **Two-Stage Framework** instantiated with
  - Milne-Witten as Top-k Retrieval
  - Weights are the average between Milne-Witten and DeepWalk
- ▷ More experiments in the paper (first known comparison among *more than 15 methods!*)



# Experiments

## Extrinsic Evaluation

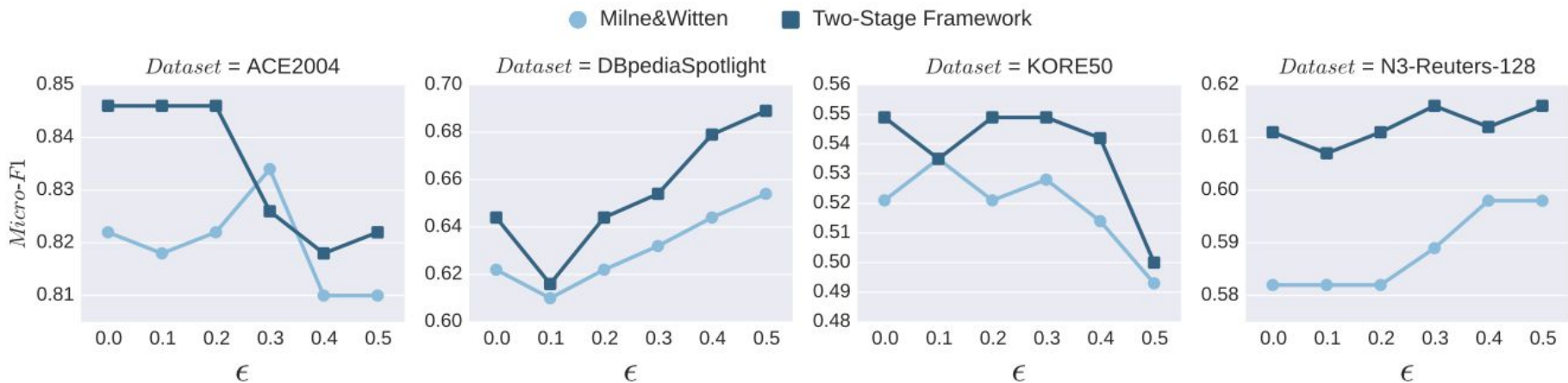


### ▷ Domain of *Entity Linking*

- Linking short but meaningful sequence of words with proper Wikipedia Entities

### ▷ Entity Linker used for experiments: **TAG-ME**

- We replaced the relatedness method used in TagMe (e.g. [Milne-Witten](#)) with our **Two-Stage Framework**



- ### ▷ Our relatedness measure not only improves TagMe, but also makes it more insensitive to choices of the $\epsilon$ -parameter in TagMe

# Experiments

## Optimizations & Efficiency

- ▷ **Top-k preprocessing** of Milne&Witten on the entities' out-neighbors
- ▷ **Compression** of
  - *Wikipedia Graph* with Webgraph (Boldi, WWW '04)
  - *DeepWalk embeddings* with FEL (Blanco, WSDM '15)

	<b>Uncompressed</b>	<b>Compressed</b>	
<b>Average Time</b>	0.5 ms	3 ms	6x slower
<b>Space</b>	5 GB	445 MB	10x space-saving!

Our framework fits in few hundred of MB and the computation of the relatedness is still sufficiently fast at query time!

# 2

# Algorithms for Entity and Fact Salience

Document Aboutness via Sophisticated Syntactic  
and Semantic Features

Marco Ponza, Paolo Ferragina, and Francesco Piccinno



SWAT: A System for Detecting  
Salient Wikipedia Entities in Texts

Marco Ponza, Paolo Ferragina, and Francesco Piccinno



Facts That Matter

Marco Ponza, Luciano Del Corro, and Gerhard Weikum



# Introduction

## Automatic Document Summarization



- ▷ Succinct representation of the **Document's Subject Matter** (Bruza, AIR '96)
  
- ▷ Condensing **salient information** from an input text into a **summary**
  - Enable fast and accurate document search
  - Help reader to identify relevant topics

(Hasan, ACL'14)

# Introduction

## Automatic Document Summarization



- ▷ Succinct representation of the **Document's Subject Matter** (Bruza, AIR '96)
- ▷ Condensing **salient information** from an input text into a **summary**



Input Document



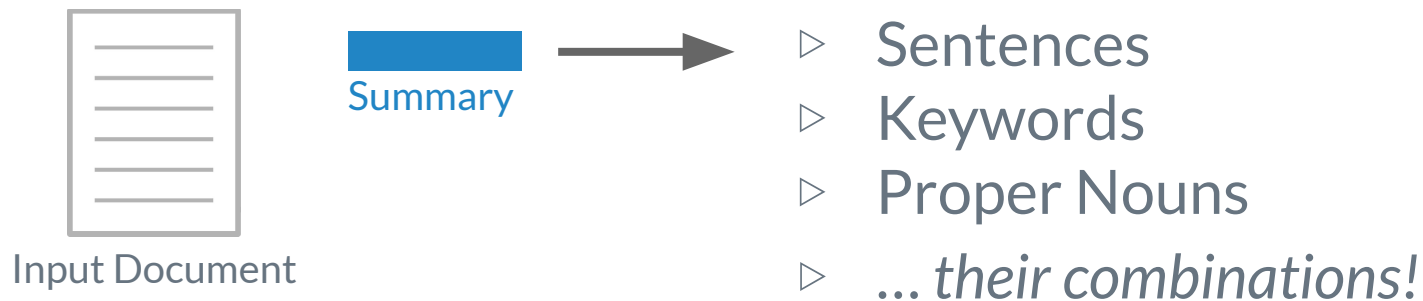
Summary

# Introduction

## Automatic Document Summarization



- ▷ Succinct representation of the **Document's Subject Matter** (Bruza, AIR '96)
- ▷ Condensing **salient information** from an input text into a **summary**



Input Document

Summary

- ▷ Sentences
- ▷ Keywords
- ▷ Proper Nouns
- ▷ ... *their combinations!*

# Introduction

## Automatic Document Summarization



- ▷ Succinct representation of the **Document's Subject Matter** (Bruza, AIR '96)
- ▷ Condensing **salient information** from an input text into a **summary**

### Our Context



Input Document



Summary

We attack the *summarization* problem from *two different points of view*

- ▷ Salient Wikipedia Entities
- ▷ Salient Open Facts

# 2.1

# Algorithms for Entity Salience

Document Aboutness via Sophisticated Syntactic and Semantic Features

Marco Ponza, Paolo Ferragina and Francesco Piccinno

SWAT: A System for Detecting Salient Wikipedia Entities in Texts

Marco Ponza, Paolo Ferragina and Francesco Piccinno





# Entity Saliience

Summarization via Salient Wikipedia Entities



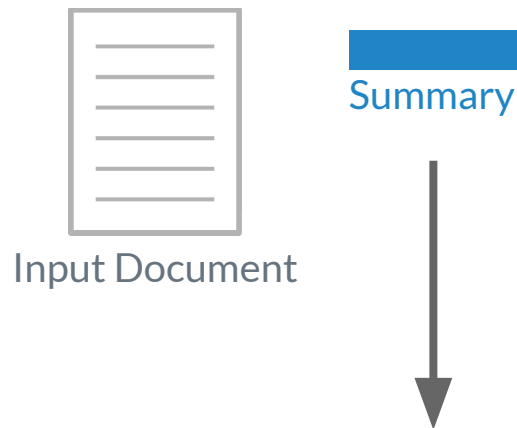
Input Document

Summary



# Entity Salience

Summarization via Salient Wikipedia Entities



Set of Salient Wikipedia Entities



# Entity Saliency

## Contributions: Our Solution



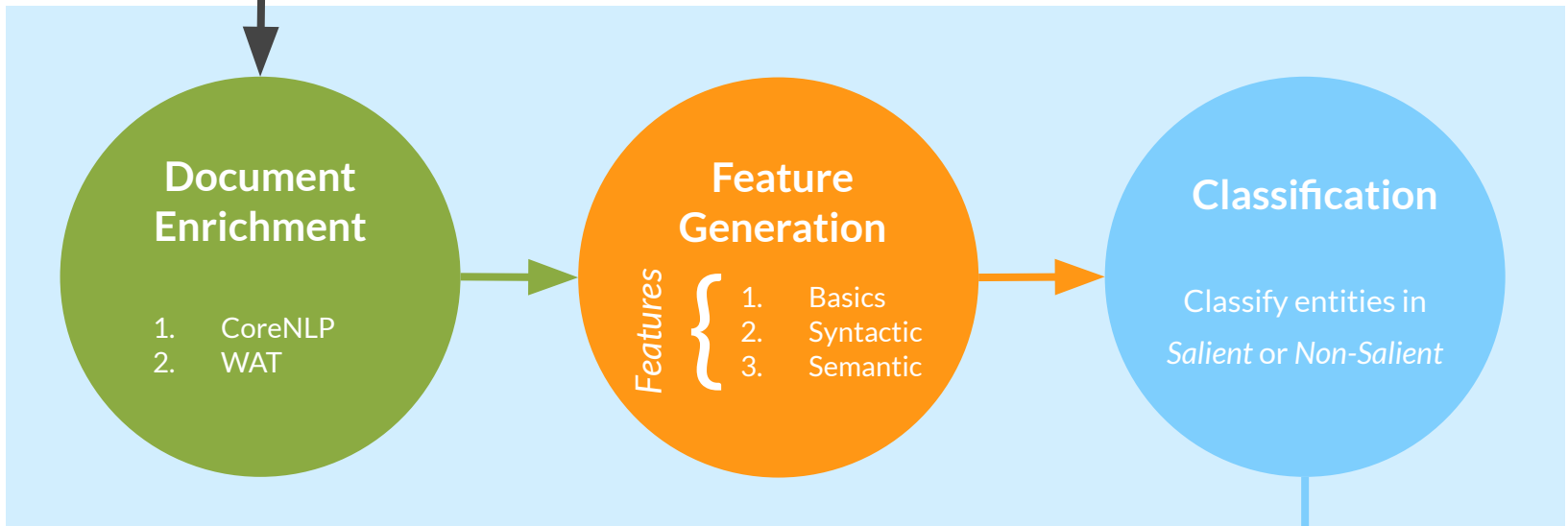
- ▷ **Three-Stage System** for Entity Saliency Extraction
- ▷ In-Depth Feature Engineering:
  - **Syntactic:**
    - Sentence Ranking
    - Dependency Trees
  - **Semantic:**
    - Entity Annotations
    - Relatedness Graph
- ▷ **Improves** current solutions
  - From +1.9% up to +14%
- ▷ The **first publicly** available **API**

# Entity Saliency

## General Structure

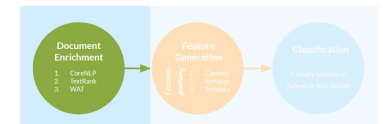


*Input Document*



*Salient Entities*





# Entity Salience

## Three-Stage System

The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

### POLITICAL ACTION; Decisions on the Horizon

By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators Barack Obama and Hillary Rodham Clinton anytime soon, but stay tuned.

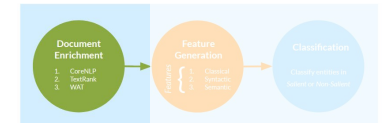
At least that is the word from their associates. Mr. Obama, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after President Bush's State of the Union address on Jan. 23. As he walked out of the Capitol on a recent afternoon, Mr. Obama only smiled when asked about his timing. Then, he rushed to change the subject.

Initially, Mr. Obama said he intended to announce his decision after returning from a holiday vacation in Hawaii, where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to make up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS

# 1. Document Enrichment

▶ CoreNLP (Manning, ACL '14)



# Entity Salience

## Three-Stage System

The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

### POLITICAL ACTION; Decisions on the Horizon

By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators Barack Obama and Hillary Rodham Clinton anytime soon, but stay tuned.

At least that is the word from their associates. Mr. Obama, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after President Bush's State of the Union address on Jan. 23. As he walked out of the Capitol on a recent afternoon, Mr. Obama only smiled when asked about his timing. Then, he rushed to change the subject.

Initially, Mr. Obama said he intended to announce his decision after returning from a holiday vacation in Hawaii, where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to make up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS

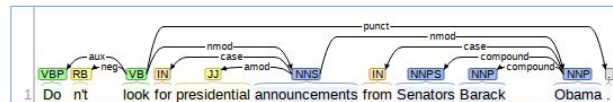
#### Part-of-Speech:

VBVP RB VB IN JJ NNS IN NNPS NNP NNP  
1 Do nt look for presidential announcements from Senators Barack Obama .

#### Named Entity Recognition:

PERSON  
1 Do nt look for presidential announcements from Senators Barack Obama .

#### Basic Dependencies:



# 1. Document Enrichment

▷ CoreNLP (Manning, ACL '14)

## Module

Sentence Splitting

Tokenization

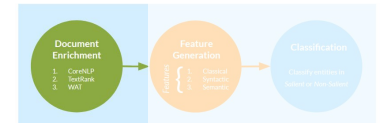
POS-Tagging

Named Entity Recognition

Dependency Parsing

Coreference

Images via  
<http://corenlp.run>



# Entity Salience

## Three-Stage System

The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

### POLITICAL ACTION; Decisions on the Horizon

By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators **Barack Obama** and **Hillary Rodham Clinton** anytime soon, but stay tuned.

At least that is the word from their associates. **Mr. Obama**, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after **President Bush's** State of the Union address on Jan. 23. As he walked out of the **Capitol** on a recent afternoon, **Mr. Obama** only smiled when asked about his timing. Then, he rushed to change the subject.

Initially, **Mr. Obama** said he intended to announce his decision after returning from a holiday vacation in **Hawaii** where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to make up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS



# 1. Document Enrichment

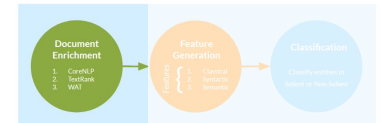
▷ **CoreNLP** (Manning, ACL '14)

Named Entities + Proper/Common Nouns



▷ **WAT** (Piccinno, SIGIR '14)

- Annotates them with Wikipedia Entities



# Entity Salience

## Three-Stage System

The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

### POLITICAL ACTION; Decisions on the Horizon

By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators **Barack Obama** and **Hillary Rodham Clinton** anytime soon, but stay tuned.

At least that is the word from their associates. **Mr. Obama**, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after **President Bush's** State of the Union address on Jan. 23. As he walked out of the **Capitol** on a recent afternoon, **Mr. Obama** only smiled when asked about his timing. Then, he rushed to change the subject.

Initially, **Mr. Obama** said he intended to announce his decision after returning from a holiday vacation in **Hawaii** where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to make up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS



# 1. Document Enrichment

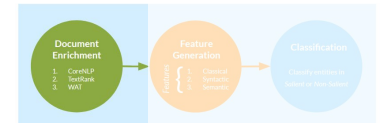
▷ **CoreNLP** (Manning, ACL '14)

Named Entities + Proper/Common Nouns

▷ **WAT** (Piccinno, SIGIR '14)

- Annotates them with Wikipedia Entities





# Entity Salience

## Three-Stage System

The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

### POLITICAL ACTION; Decisions on the Horizon

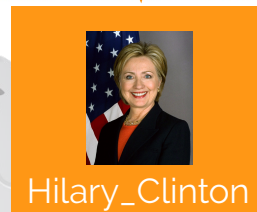
By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators **Barack Obama** and **Hillary Rodham Clinton** any time soon, but stay tuned.

At least that is the word from their associates. **Mr. Obama**, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after **President Bush's** State of the Union address on Jan. 23. As he walked out of the **Capitol** on a recent afternoon, **Mr. Obama** only smiled when asked about his timing. Then, he rushed to change the subject.

Initially, **Mr. Obama** said he intended to announce his decision after returning from a holiday vacation in **Hawaii** where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to wake up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS



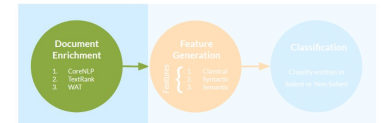
# 1. Document Enrichment

▷ **CoreNLP** (Manning, ACL '14)

Named Entities + Proper/Common Nouns

▷ **WAT** (Piccinno, SIGIR '14)

- Annotates them with Wikipedia Entities



# Entity Salience Three-Stage System

The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

## POLITICAL ACTION; Decisions on the Horizon

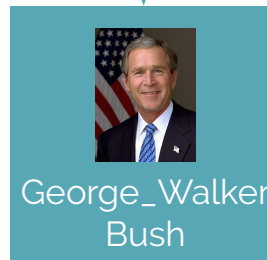
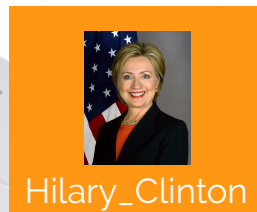
By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators **Barack Obama** and **Hillary Rodham Clinton** anytime soon, but stay tuned.

At least that is the word from their associates. **Mr. Obama**, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after **President Bush's** State of the Union address on Jan. 23. As he walked out of the **Capitol** on a recent afternoon, **Mr. Obama** only smiled when asked about his timing. Then he rushed to change the subject.

Initially, **Mr. Obama** said he intended to announce his decision after returning from a holiday vacation in **Hawaii** where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to make up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS



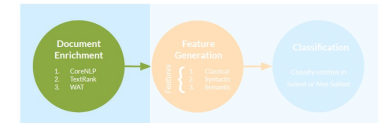
# 1. Document Enrichment

▷ **CoreNLP** (Manning, ACL '14)

Named Entities + Proper/Common Nouns

▷ **WAT** (Piccinno, SIGIR '14)

- Annotates them with Wikipedia Entities



# Entity Salience

## Three-Stage System

The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

### POLITICAL ACTION; Decisions on the Horizon

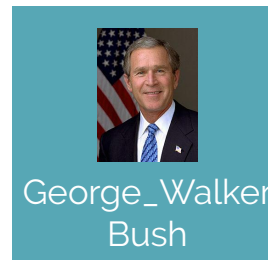
By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators **Barack Obama** and **Hillary Rodham Clinton** anytime soon, but stay tuned.

At least that is the word from their associates. **Mr. Obama**, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after **President Bush's** State of the Union address on Jan. 23. As he walked out of the **Capitol** on a recent afternoon, **Mr. Obama** only smiled when asked about his timing. Then, he rushed to change the subject.

Initially, **Mr. Obama** said he intended to announce his decision after returning from a holiday vacation in **Hawaii** where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to make up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS



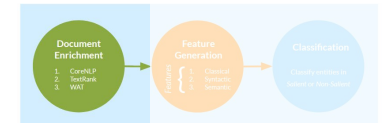
# 1. Document Enrichment

▷ CoreNLP (Manning, ACL '14)

Named Entities + Proper/Common Nouns

▷ WAT (Piccinno, SIGIR '14)

- Annotates them with Wikipedia Entities



# Entity Salience

## Three-Stage System

The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

### POLITICAL ACTION; Decisions on the Horizon

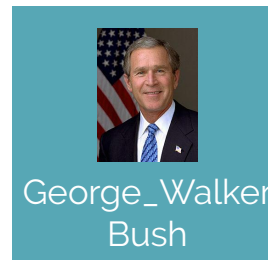
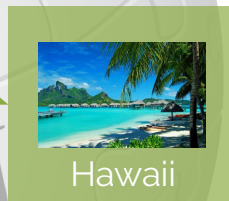
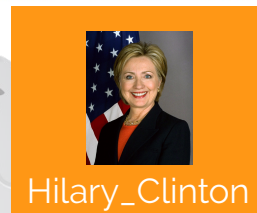
By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators **Barack Obama** and **Hillary Rodham Clinton** anytime soon, but stay tuned.

At least that is the word from their associates. **Mr. Obama**, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after **President Bush's** State of the Union address on Jan. 23. As he walked out of the **Capitol** on a recent afternoon, **Mr. Obama** only smiled when asked about his timing. Then, he rushed to change the subject.

Initially, **Mr. Obama** said he intended to announce his decision after returning from a holiday vacation in **Hawaii** where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to make up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS



# 1. Document Enrichment

▷ CoreNLP (Manning, ACL '14)

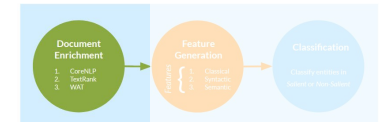
Named Entities + Proper/Common Nouns

▷ WAT (Piccinno, SIGIR '14)

- Annotates them with Wikipedia Entities

# Entity Salience

## Three-Stage System



The New York Times

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION

### POLITICAL ACTION; Decisions on the Horizon

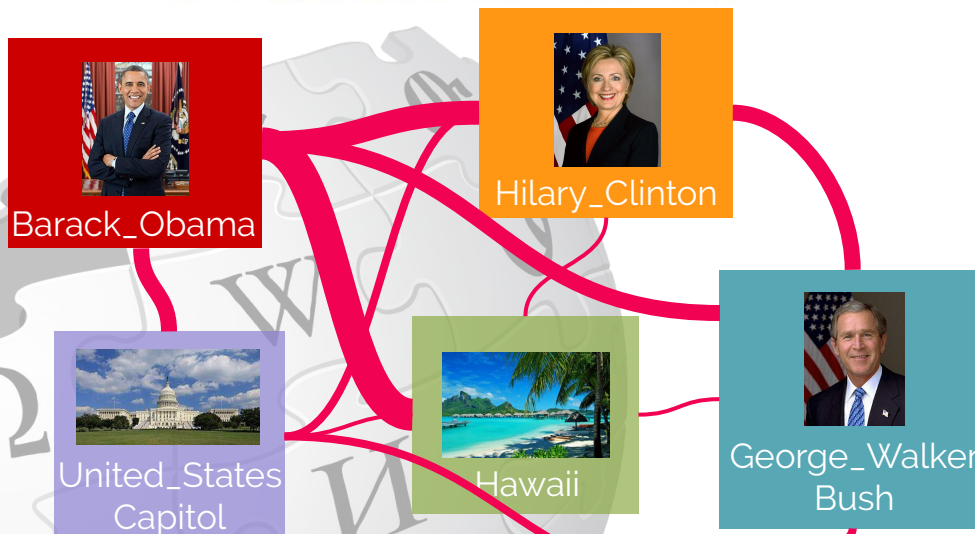
By JEFF ZELENY and PATRICK HEALY  
Published: January 9, 2007

Don't look for presidential announcements from Senators **Barack Obama** and **Hillary Rodham Clinton** anytime soon, but stay tuned.

At least that is the word from their associates. **Mr. Obama**, Democrat of Illinois, is not likely to say whether he intends to seek the party's presidential nomination until after **President Bush's** State of the Union address on Jan. 23. As he walked out of the **Capitol** on a recent afternoon, **Mr. Obama** only smiled when asked about his timing. Then, he rushed to change the subject.

Initially, **Mr. Obama** said he intended to announce his decision after returning from a holiday vacation in **Hawaii** where he was visiting his grandmother and other relatives. Now, several people close to the senator say, he needs a little more time to make up his mind.

- FACEBOOK
- TWITTER
- GOOGLE+
- EMAIL
- SHARE
- PRINT
- REPRINTS



# 1. Document Enrichment

▷ **CoreNLP** (Manning, ACL '14)

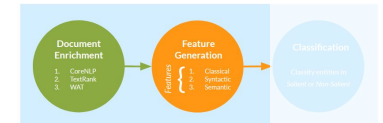
Named Entities + Proper/Common Nouns

▷ **WAT** (Piccinno, SIGIR '14)

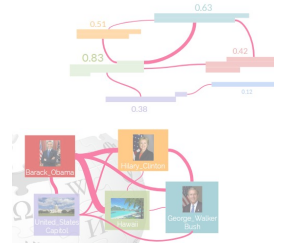
- Annotates them with Wikipedia Entities
- Relatedness Graph
  - Nodes = Entities
  - **Weights** are defined with two different relatedness scores:
    - (i) Wikipedia **Jaccard** In-Links
    - (ii) **Cosine** between Entity Embeddings

# Entity Salience

## Three-Stage System



Tokens, POS Tags, Dependency Relations, Coreference Chains, Wikipedia Entities and their Relatedness



## 2. Feature Generation

### ▷ Standard Entity Features

- Frequency
- Positions
- ...

### ▷ Syntactic Features

- Statistics on Sentence Ranks
- Frequency/Positions of Dependency Relations
- ...

### ▷ CMU-Google Features

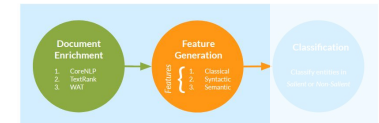
- POS-Tags, Coreference Freq.
- PageRank on a graph whose weights are based on co-occ.
- ...

### ▷ Semantic Features

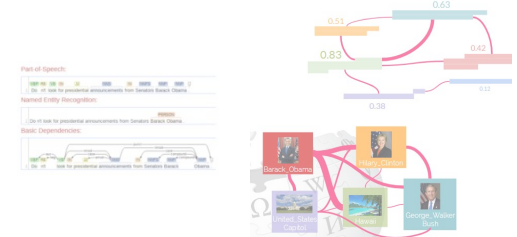
- Statistics on annotations (coherence, commonness)
- Graph Centralities on Relatedness Graph
- Relatedness over Positions
- ...

# Entity Salience

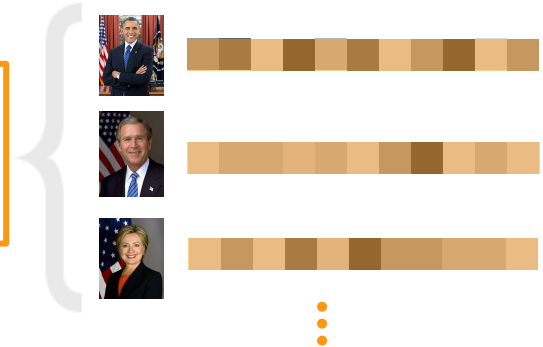
## Three-Stage System



Tokens, POS Tags, Dependency Relations, Coreference Chains, Sentence Ranks, Wikipedia Entities and their Relatedness



Entity Feature Vectors



Salient Entities



dmlc  
**XGBoost**  
(Chen, SIGKDD '16)



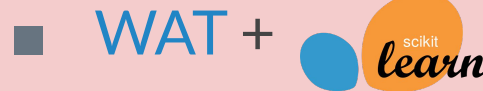
# Experiment

## Competitors & Benchmarks

- ▷ **CMU-Google System** (Dunietz, EACL '14)

### Re-implemented

- Proprietary modules substituted with open-source tools



- ▷ **SEL** (Trani, DocEng '16)


Limitations

- No comparison with CMU-Google System
- Benchmark on small dataset
- Not publicly available

Dataset (365 news, 4747 entities)

# Experiments

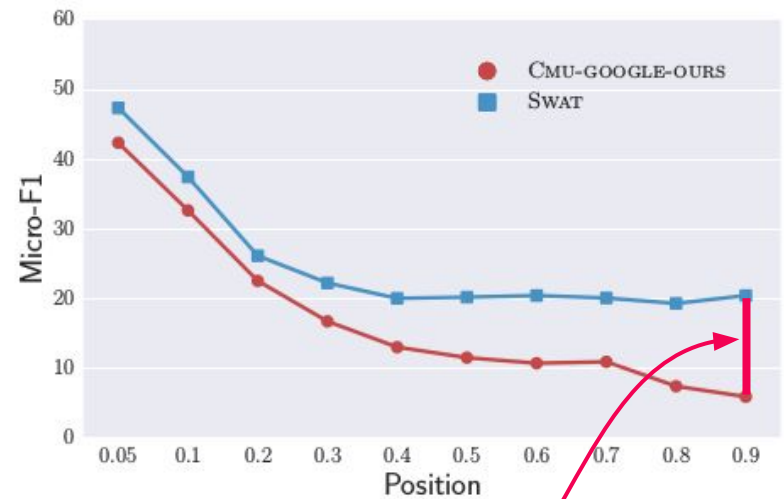
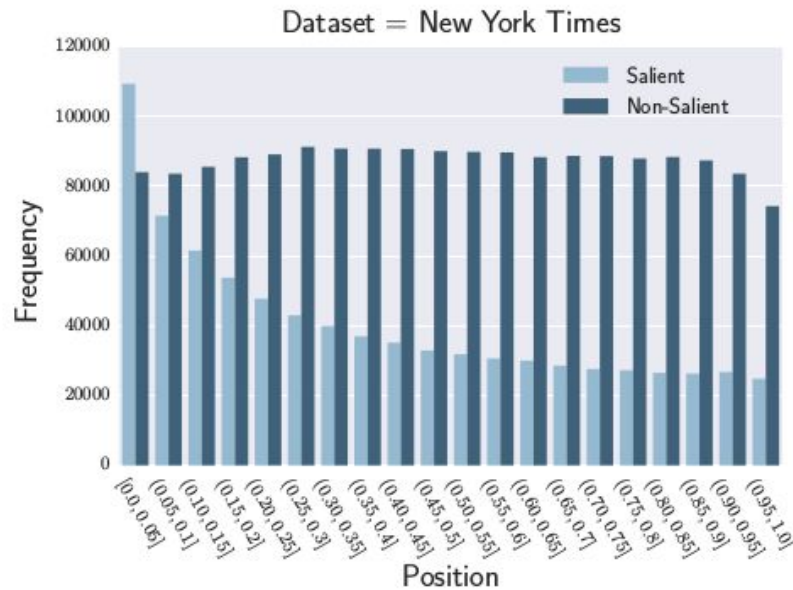
## Results

System	New York Times			Wikinews		
	Micro			Macro		
	P	R	F1	P	R	F1
CMU-Google (Dunietz, EACL '14)	60.5	63.5	61.5	-	-	-
CMU-Google-ours	58.8	62.6	60.7	42.3	61.0	46.0
SEL (Trani, DocEng '16)	-	-	-	<b>61.0</b>	50.0	52.0
 SWAT	<b>62.4</b>	<b>66.0</b>	<b>64.1</b>	57.7	<b>67.0</b>	<b>58.3</b>
	+1.9%	+2.5%	+2.6%	-3.3%	+6.0%	+6.3%

# Experiments

## Results

### Independence from position of salient entities



+14%

# 2.2

# Algorithms for Fact Salience

Facts That Matter \*

Marco Ponza, Luciano Del Corro and Gerhard Weikum

\*work done during an internship at

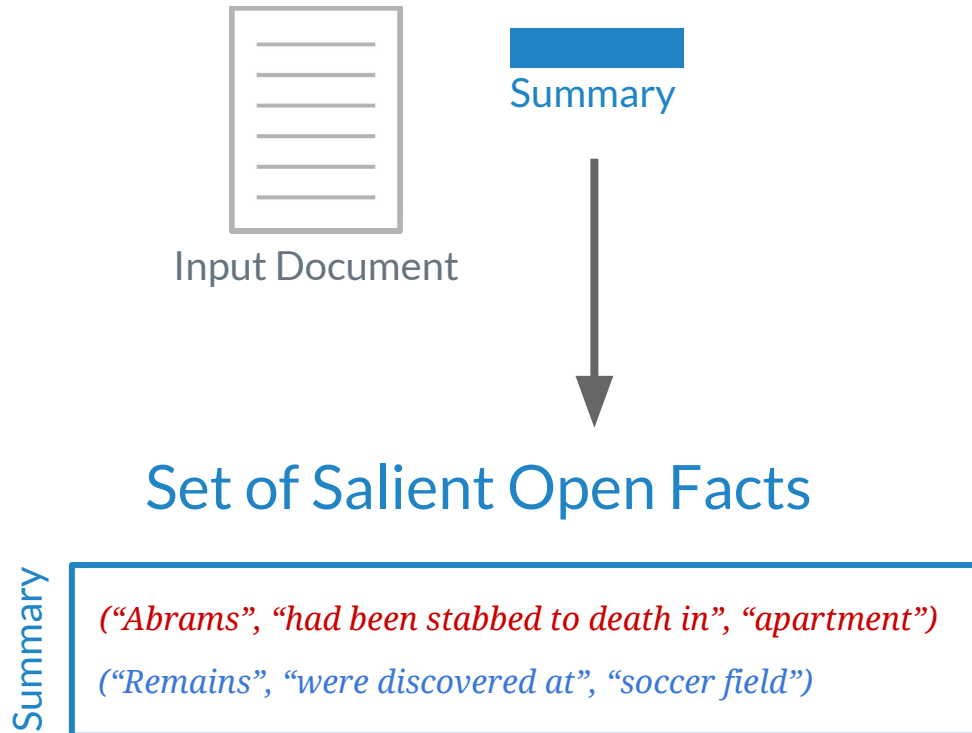


**EMNLP 2018**

Brussels

# Fact Salience

Summarization via Salient Open Facts



# Fact Salience

## Contributions

1

We introduce a **NEW** Research Task called **Fact Salience**

- Extraction of **relevant information** from an input document expressed in the **smallest number of facts**



2

Proposal of the **1st Fact Salience system**:

**SalIE**

Fully **unsupervised**

Based on **PageRank** and **Clustering**

**Public available at**

<https://github.com/mponza/SalIE>



3

Experiments show that open facts are an **effective way to compress information**

# SalIE Salient Information Extraction Overview



Input Document



Open Information Extraction  
MinIE (Gashteovski, EMNLP 2017)

Open Facts

*(“Abrams”, “was 56-years-old native of”, “Pittsburgh area”)*

*(“Abrams”, “had been stabbed to death in”, “apartment”)*

*(“Apartment”, “tending wounds at time of”, “murder”)*

*(“Cousin of husband”, “had gone into”, “business”)*

*(“Remains”, “were discovered at”, “soccer field”)*

*(“Abrams”, “got more involved in”, “real estate”)*

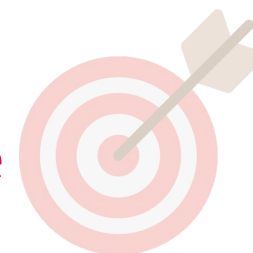


Salient Open Facts

# SalIE Salient Information Extraction

First Stage: Fact Relevance

Provide, for each open fact, a **relevance score**



*“Most relevant facts are the ones more central in the input document”*



We can use PageRank!

1. How do we define the **graph structure**?
2. How do we **weight** the **edges**?
3. How do we instantiate the **teleport vector**?



# SalIE Salient Information Extraction

## First Stage: Fact Relevance

Provide, for each open fact, a **relevance score**



Open Facts

*("Abrams", "was 56-years-old native of", "Pittsburgh area")*

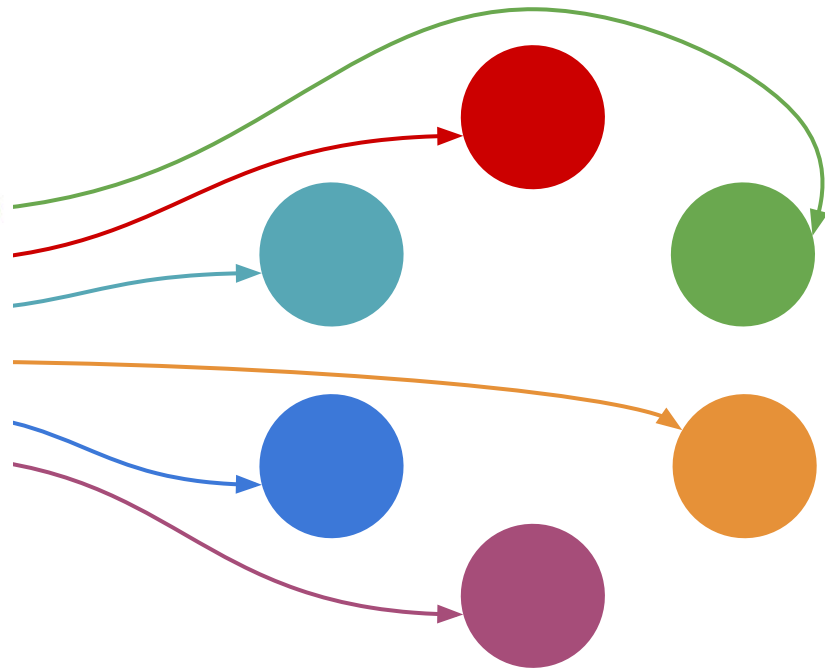
*("Abrams", "had been stabbed to death in", "apartment")*

*("Apartment", "tending wounds at time of", "murder")*

*("Cousin of husband", "had gone into", "business")*

*("Remains", "were discovered at", "soccer field")*

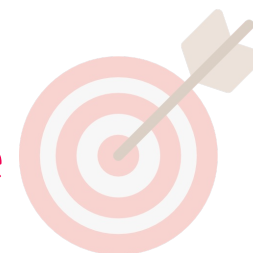
*("Abrams", "got more involved in", "real estate")*



# SalIE Salient Information Extraction

## First Stage: Fact Relevance

Provide, for each open fact, a **relevance score**



Open Facts

*(“Abrams”, “was 56-years-old native of”, “Pittsburgh area”)*

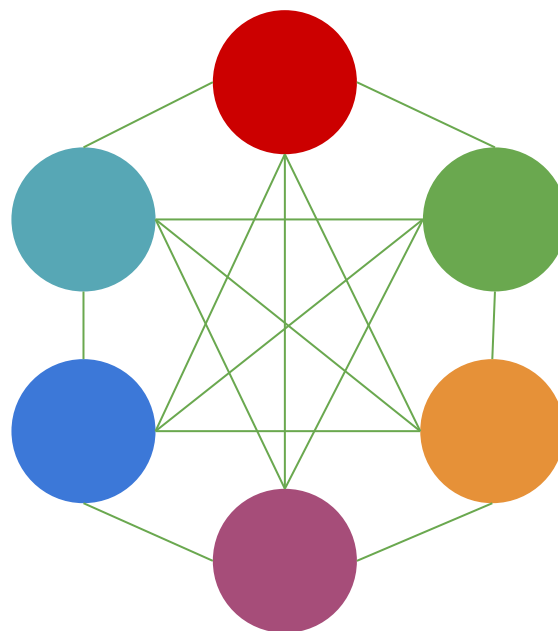
*(“Abrams”, “had been stabbed to death in”, “apartment”)*

*(“Apartment”, “tending wounds at time of”, “murder”)*

*(“Cousin of husband”, “had gone into”, “business”)*

*(“Remains”, “were discovered at”, “soccer field”)*

*(“Abrams”, “got more involved in”, “real estate”)*



1. How do we define the graph structure?  
We can grow a *fully connected graph of facts!*

# SalIE Salient Information Extraction

## First Stage: Fact Relevance

Provide, for each open fact, a **relevance score**



Open Facts

*("Abrams", "was 56-years-old native of", "Pittsburgh area")*

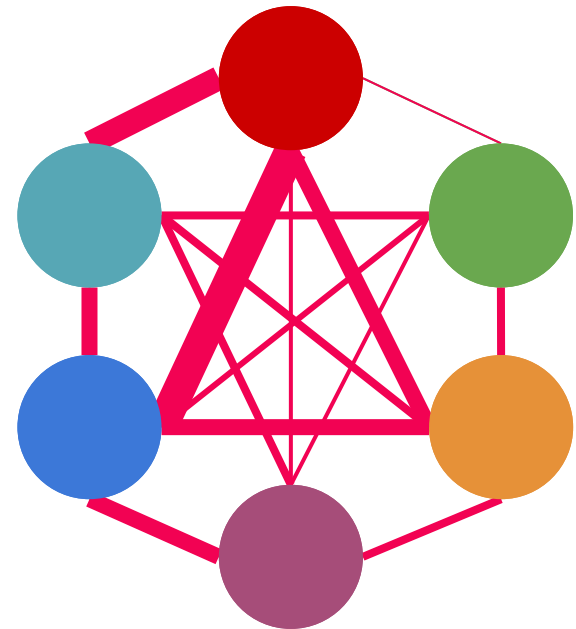
*("Abrams", "had been stabbed to death in", "apartment")*

*("Apartment", "tending wounds at time of", "murder")*

*("Cousin of husband", "had gone into", "business")*

*("Remains", "were discovered at", "soccer field")*

*("Abrams", "got more involved in", "real estate")*

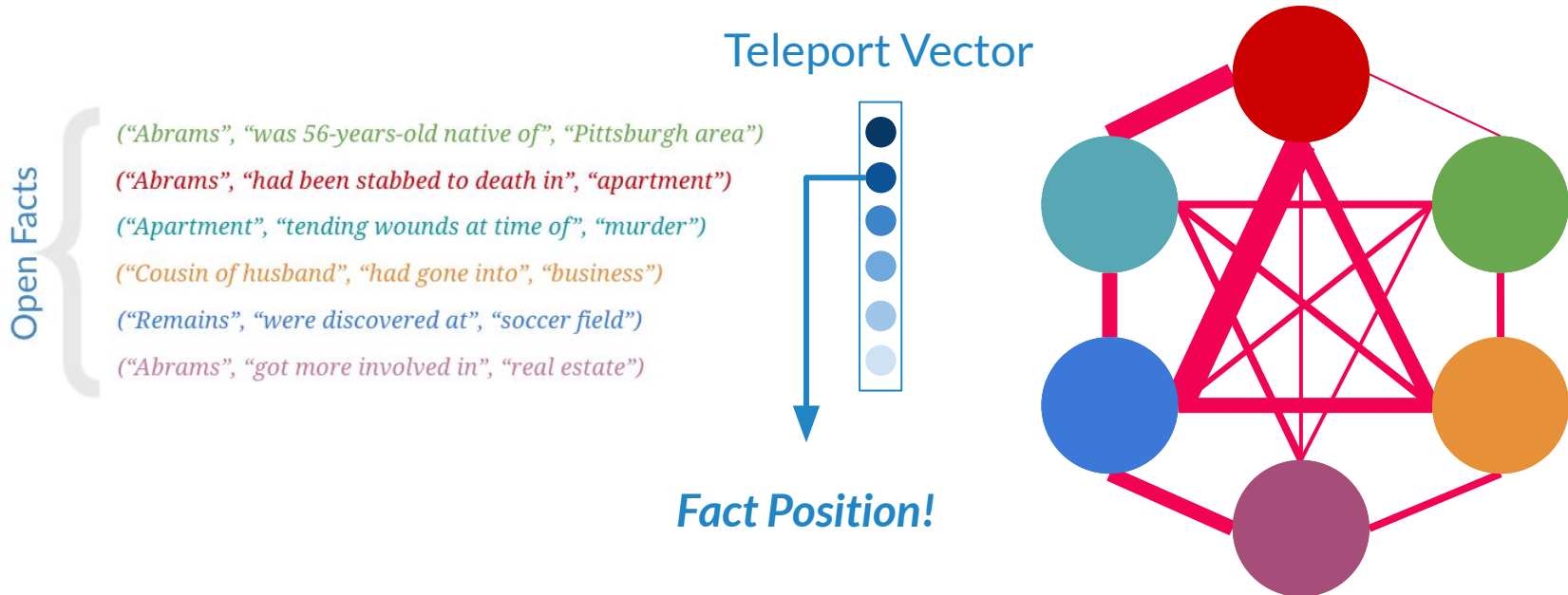


2. How do we **weight** the **edges**?  
We can use the **cosine similarity** between **facts' embeddings vectors**

# SalIE Salient Information Extraction

## First Stage: Fact Relevance

Provide, for each open fact, a **relevance score**

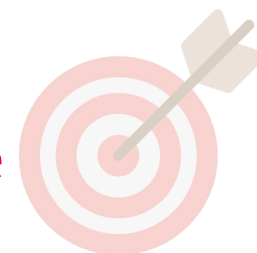


3. How do we instantiate the **teleport vector**?  
Each facts' entry in the teleport vector is scored wrt facts' **positional** information

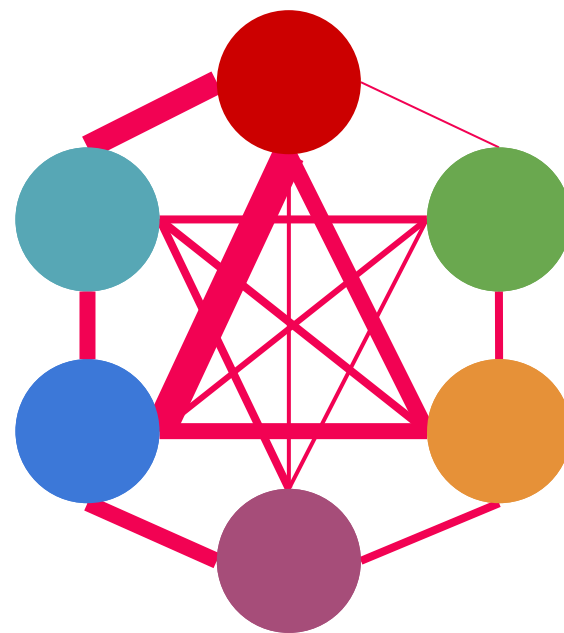
# SalIE Salient Information Extraction

## First Stage: Fact Relevance

Provide, for each open fact, a **relevance score**



- ✓ 1. Fully connected graph
- ✓ 2. Edges weighted with cosine between word embeddings
- ✓ 3. Teleport vector instantiated as a function of the position



Yes! Now we can run PageRank!

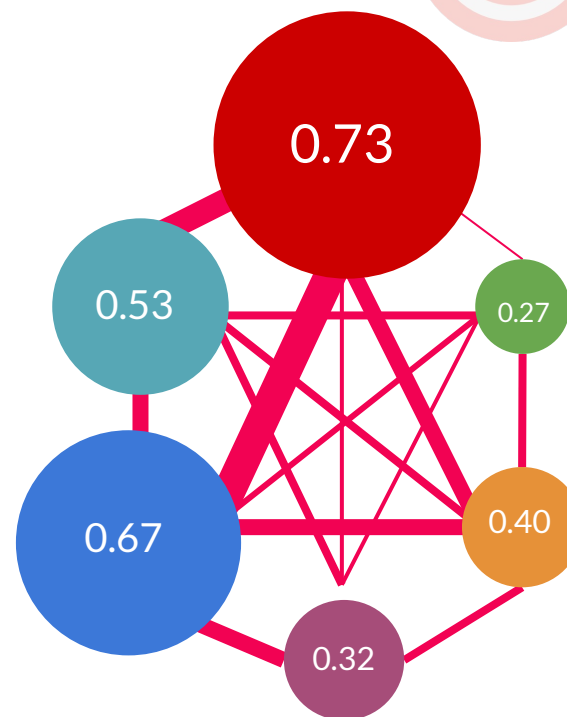
# SalIE Salient Information Extraction

First Stage: Fact Relevance

Provide, for each open fact, a **relevance score**



- 1. Fully connected graph
- 2. Edges weighted with cosine between word embeddings
- 3. Teleport vector instantiated as a function of the position



Yes! Now we can run PageRank!

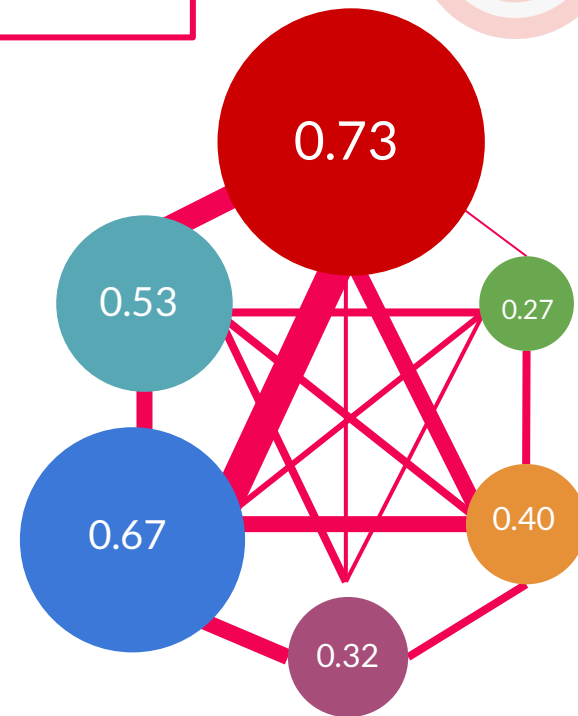
# SalIE Salient Information Extraction

## First Stage: Fact Relevance

Provide, for each open fact, a **relevance score**



<i>("Abrams", "was 56-years-old native of", "Pittsburgh area")</i>	0.27
<i>("Abrams", "had been stabbed to death in", "apartment")</i>	0.73
<i>("Apartment", "tending wounds at time of", "murder")</i>	0.53
<i>("Cousin of husband", "had gone into", "business")</i>	0.40
<i>("Remains", "were discovered at", "soccer field")</i>	0.67
<i>("Abrams", "got more involved in", "real estate")</i>	0.32



# SalIE Salient Information Extraction

Second Stage: Fact Diversification

Maximize the information in the  
*smallest number of facts*



- ▷ Salient open facts should provide a wide spectrum of the information in the document content



Salient Fact **Diversification** with **Clustering**



# SalIE Salient Information Extraction

## Second Stage: Fact Diversification

Maximize the information in the smallest number of facts



*("Abrams", "was 56-years-old native of", "Pittsburgh area")* 0.27  
*("Abrams", "had been stabbed to death in", "apartment")* 0.73  
*("Apartment", "tending wounds at time of", "murder")* 0.53  
*("Cousin of husband", "had gone into", "business")* 0.40  
*("Remains", "were discovered beside warehouse at edge of", "cinder-topped soccer field on outskirts of Panama City")* 0.67  
*("Abrams", "got more involved in", "real estate")* 0.32

Ranked Facts  
(from the First Stage)



### Clustering Rule

Facts are clustered together wrt their subject

0.73 - (*"Abrams", "had been stabbed to death in", "apartment"*)

0.32 - (*"Abrams", "got more involved in", "real estate"*)

0.27 - (*"Abrams", "was 56-years-old native of", "Pittsburgh area"*)

0.67 - (*"Remains", "were discovered at", "soccer field"*)

0.53 - (*"Apartment", "tending wounds at time of", "murder"*)

0.40 - (*"Cousin of husband", "had gone into", "business"*)

# SalIE Salient Information Extraction

## Second Stage: Fact Diversification

Maximize the information in the  
smallest number of facts



- 0.73 - ("Abrams", "had been stabbed to death in", "apartment")
- 0.32 - ("Abrams", "got more involved in", "real estate")
- 0.27 - ("Abrams", "was 56-years-old native of", "Pittsburgh area")
- 0.67 - ("Remains", "were discovered beside warehouse at edge of", "cinder-topped soccer field on outskirts of Panama City")
- 0.53 - ("Apartment", "tending wounds at time of", "murder")
- 0.40 - ("Cousin of husband", "had gone into", "business")



### Diversification Rule

From each cluster, select the  
fact with the strongest  
PageRank score

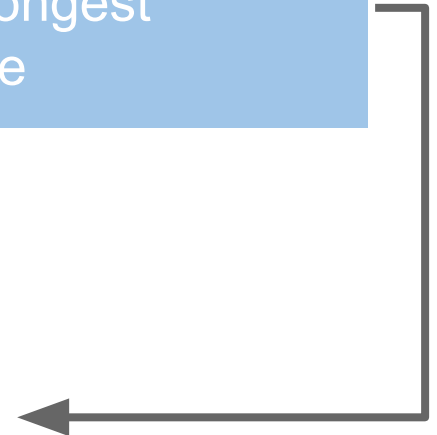
Final Output

("Abrams", "had been stabbed to death in", "apartment")

("Remains", "were discovered at", "soccer field")

("Apartment", "tending wounds at time of", "murder")

("Cousin of husband", "had gone into", "business")



# Experiments

## Evaluation

- ▷ **Dataset: New York Times**
  - ~4K news each one with a manually annotated human summary
- ▷ **Methodology:** Evaluate how **salient** are **top-k facts**
  - We evaluate 5 configurations: From top-1 to top-5 facts
- ▷ **Metrics:** How close are top-k facts to a human summary?
  - **ROUGE** (Lin et al. 2004) metrics
    - % of unigrams, bigrams, subsequences between **generated summary** and ground-truth


# Experiments

## Baselines

- ▷ **Position** Baseline
  - Returns facts wrt their **order** in the document
  - Standard for saliency and summarization tasks
- ▷ **TextRank** (Mihalcea, EMNLP'04)
  - Graph-based summarizer based on token-overlap between sentences
  - Re-implemented to work at fact-level
- ▷ **Berkeley** (Durrett, ACL'16)
  - **Supervised** summarizer based on handcrafted features and SVM

# Experiments

## Results

Method	ROUGE-1					ROUGE-L				
	1	2	3	4	5	1	2	3	4	5
Position	13.9	20.4	24.8	27.8	29.7	12.8	18.1	21.8	24.4	26.0
TextRank	15.2	21.5	24.5	26.1	26.8	13.0	17.5	19.8	21.3	22.0
Berkeley	8.5	18.0	25.4	<b>30.4</b>	<b>34.1</b>	8.00	16.3	22.5	<b>26.7</b>	<b>29.7</b>
 SalIE	17.1	24.2	<b>28.0</b>	30.0	30.9	<b>15.3</b>	<b>21.2</b>	<b>24.3</b>	26.0	26.8
	<b>+1.9</b>	<b>+2.7</b>	<b>+2.6</b>	<b>+1.9</b>	<b>+0.5</b>	<b>+2.3</b>	<b>+3.1</b>	<b>+1.9</b>	<b>+1.5</b>	<b>+0.2</b>

- ▷ General **improvements** over all metrics
- ▷ **Facts** are an **effective** way to **compress** information!

# 3

# Algorithms for Expert Finding

WISER: A Semantic Approach for Expert Finding  
in Academia based on Entity Linking  
Paolo Cifariello, Paolo Ferragina, and Marco Ponzà



*Information Systems 2019*

# Expert Finding

- ▷ **Searching** for **experts** with respect to an input topic



- Extremely challenging task: *Who is an expert?*

The notion of **expertise** is **hard** to **formalize** as well as to be **modeled** (Balog, FTIR'12)



*...so difficult that literature refers to **expertise** as “**tacit knowledge**”!*

- **Expertise** is actually carried by people in their **minds**
- **Machines** have only one way to **access** to people **expertise**



**Artifacts** (e.g., papers, emails, ...) people write to share their expertise!

# Experiments

## Contributions

▷ New Expert Finding system



- Fully **unsupervised**



- Jointly combines classical retrieval techniques with the **Wikipedia KG** via **Entity Linking**



Indexing



Every **authors' profile** is modeled through a small **Wikipedia graph**...



Query Time



...used to design **new profile-centric scoring strategies** for the **retrieval of experts!**





# Wikipedia Expertise Ranking

## Indexing





# Wikipedia Expertise Ranking

## Indexing



Authors, Documents



Documents indexed  
with Elasticsearch





# Wikipedia Expertise Ranking

## Indexing



Documents indexed  
with Elasticsearch



Indexing of pairs (Author, DocIDs)



# Wikipedia Expertise Ranking

## Indexing



Authors, Documents



Documents indexed  
with Elasticsearch



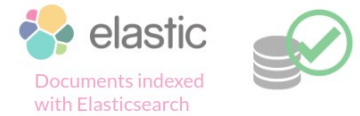
Pairs of (Author, DocIDs)



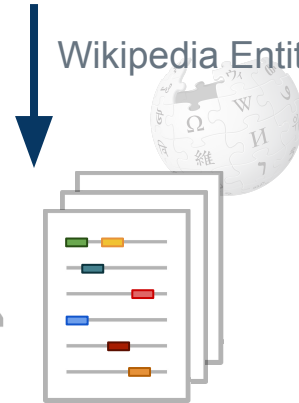


# Wikipedia Expertise Ranking

## Indexing



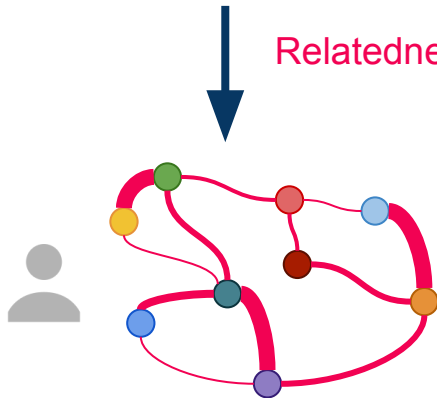
Wikipedia Entities



(Ponza, CIKM'17)



Relatedness Scores

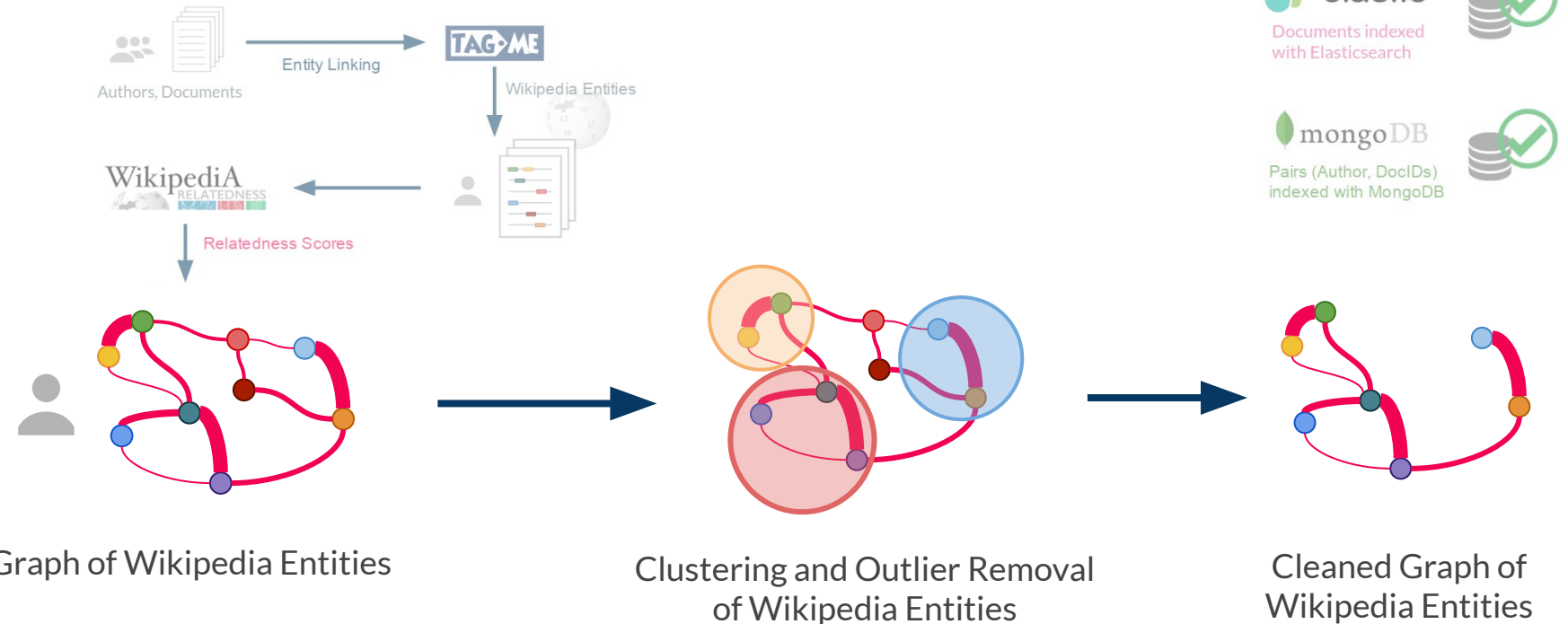


Graph of Wikipedia Entities



# Wikipedia Expertise Ranking

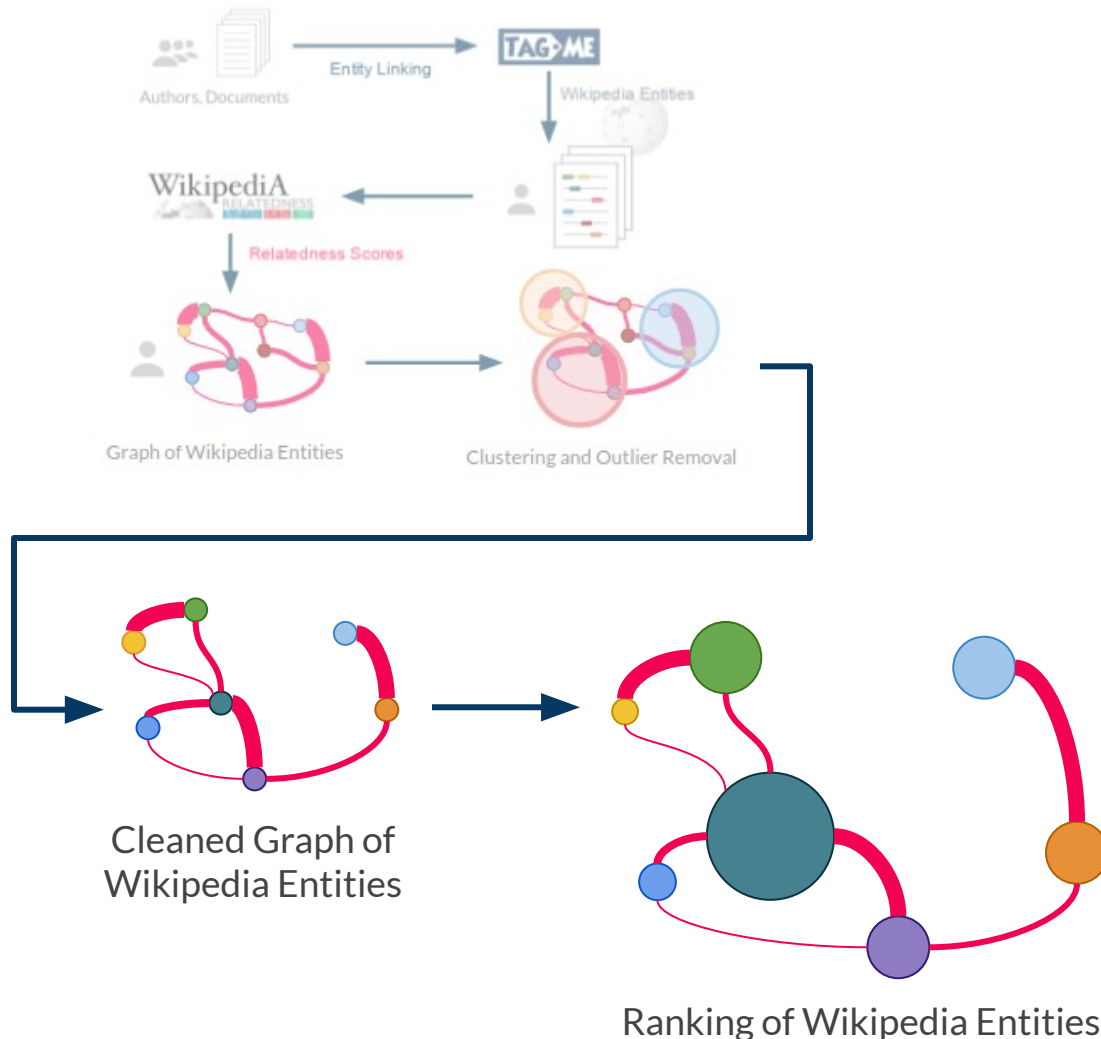
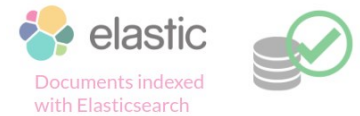
## Indexing



- ▷ **HDBSCAN** Algorithm (McInnes, IEEE'17)
- ▷ Conservative Approach:  
≥20% of nodes marked as outliers  
implies no cleaning



# Wikipedia Expertise Ranking Indexing

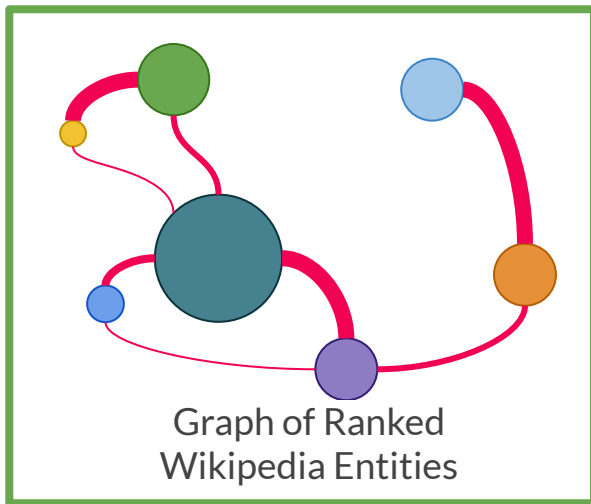


- ▷ PageRank Algorithm
- ▷ Teleport vector instantiated by taking into account the frequency of an entity annotated in the documents of the authors



# Wikipedia Expertise Ranking

## Indexing



Documents indexed  
with Elasticsearch



Pairs of (Author, DocIDs)







# Wikipedia Expertise Ranking

## Indexing

Indexing Completed!



elastic

Documents indexed  
with Elasticsearch



mongoDB

Pairs of (Author, DocIDs)

+

Pairs of (Author, Graph of Ranked  
Wikipedia Entities)





# Wikipedia Expertise Ranking

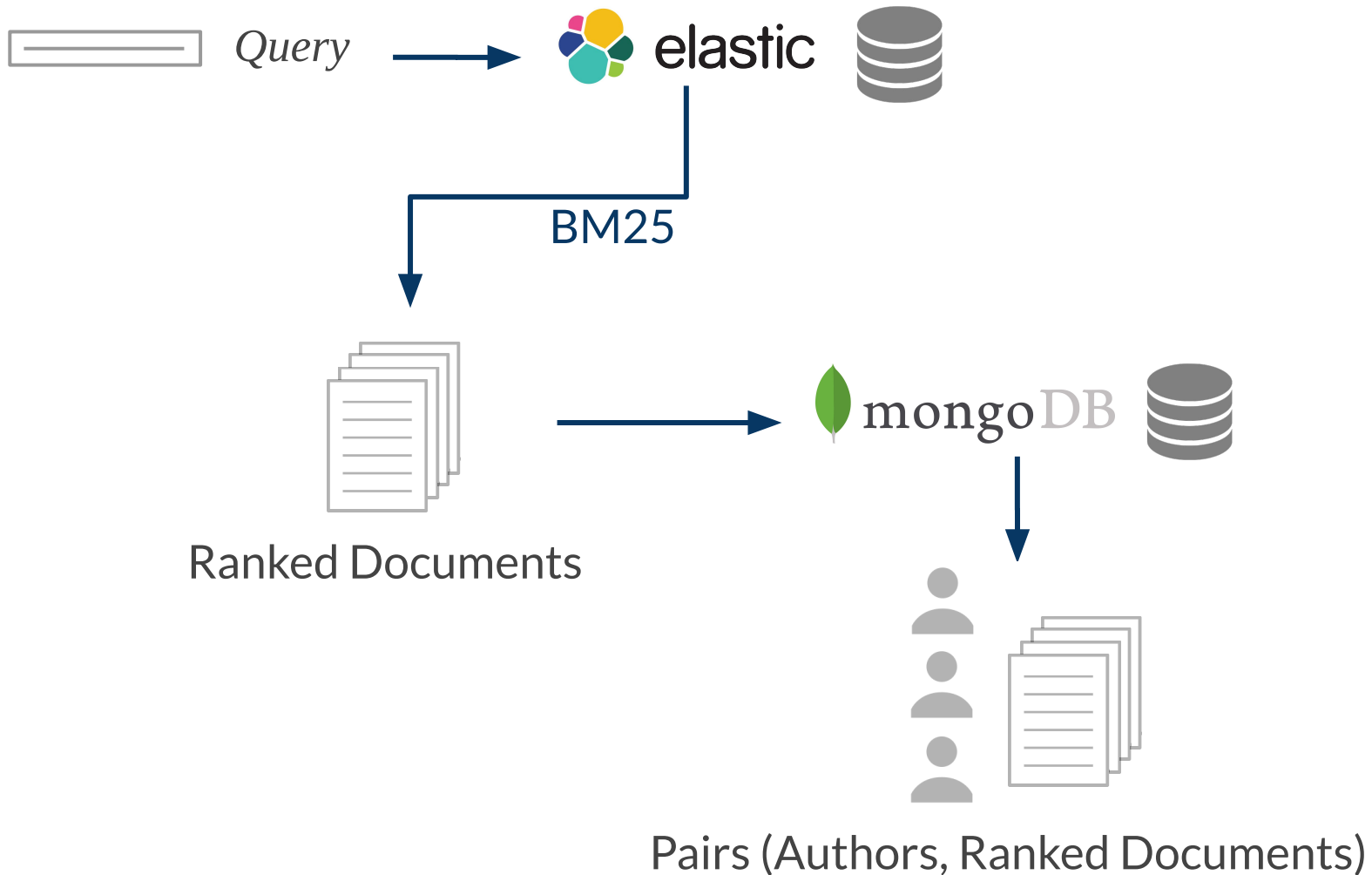
## Query Time: Two Strategies

- ▷ Jointly combine **two** different **Authors' Scoring Strategies**
  - **Document-Centric**
    1. Retrieve relevant documents
    2. Score each author wrt documents' rank (BM25)
  - **Profile-Centric**
    1. Retrieve relevant authors (wrt query Wikipedia entities)
    2. Score each author wrt entities relevance



# Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy





# Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy



1

2

3

4

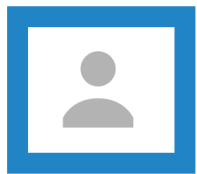
5





# Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy

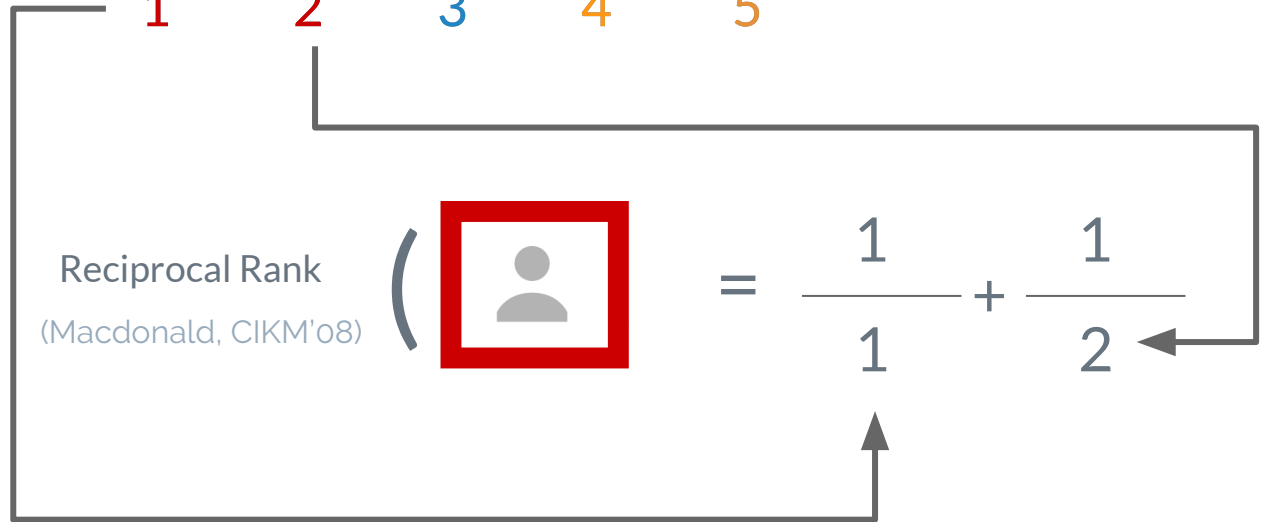


1 2 3 4 5

Reciprocal Rank  
(Macdonald, CIKM'08)



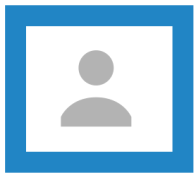
$$= \frac{1}{1} + \frac{1}{2}$$





# Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy



1 2 3 4 5

Reciprocal Rank  
(Macdonald, CIKM'08)

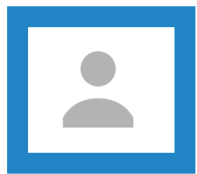


$$= \underbrace{\frac{1}{1} + \frac{1}{2}}_{1.5}$$



# Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy



1

2

3

4

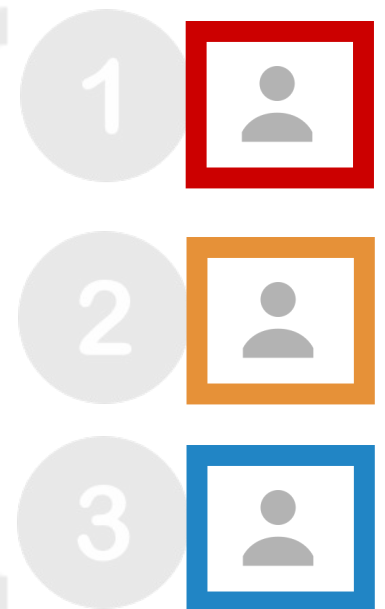
5

Reciprocal Rank (  ) = 1.5

Reciprocal Rank (  ) = 0.3

Reciprocal Rank (  ) = 0.4

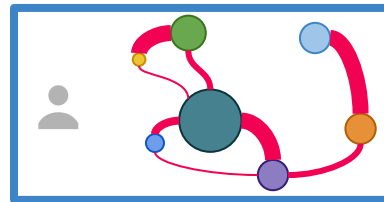
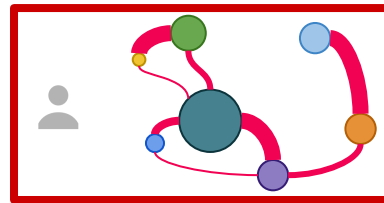
Final (Document-Centric) Ranking



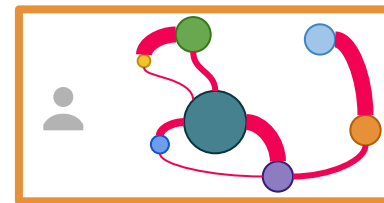


# Wikipedia Expertise Ranking

Query Time: Profile-Centric Strategy



⋮



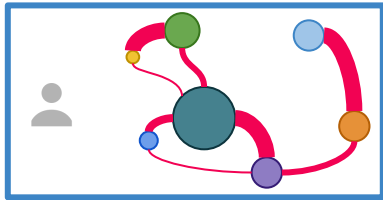
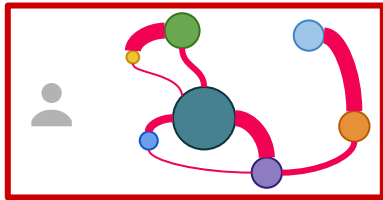
Candidate experts and their Wikipedia-based profiles matching the query's entities



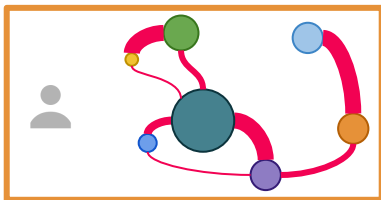


# Wikipedia Expertise Ranking

## Query Time: Profile-Centric Strategy



⋮



REC-IAF



- ▷ Score each author's entity (matched in the input query)
- ▷ Combination of multiple scores of the **entity**:
  - Document Frequency
  - TagMe Confidence in Entity Linking
  - Inverse document frequency
  - PageRank in the author's profile



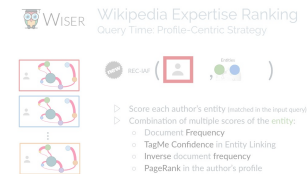
# Wikipedia Expertise Ranking

## Query Time: Data Fusion

- ▶ We have two different rankings
  - Document-Centric Ranking



- Profile-Centric Ranking



Final ranking of experts is given by the **Reciprocal Rank** (Macdonald, CIKM'08) between the **Product** of these two ranking scores

$$\text{Final Score} \left( \text{Expert} \right) = \frac{1}{\text{Doc-Cent Rank}(\text{Expert})} \cdot \frac{1}{\text{Prof-Cent Rank}(\text{Expert})}$$


# Experiments

## Benchmark

- ▷ **TU Dataset** (Berendsen, DBWIR'13)
  - ~31K documents (*largest available*)
  - ~1K researchers
  - ~1K test queries
  - Human-assessed Ground-Truth
  
- ▷ **Other systems**
  - **JM Model** (Balog, SIGIR'06)
    - Based on Frequency statistics between (Author, Keywords)
  - **Log-Linear** (Van Gysel, WWW'16)
    - Based Deep Learning (each author's profile is represented with an embedding vector)
  - **Ensemble**
    - Product Reciprocal Rank between **JM Model** and **Log-Linear**

# Experiments

## Results

Method	MAP	MRR	P@5	P@10	NDCG@100
JM Model	0.253	0.302	0.108	0.081	0.394
Log-Linear	0.287	0.363	0.134	0.092	0.425
Ensemble	0.331	0.402	0.156	<b>0.105</b>	0.477
 WISER	<b>0.385</b>	<b>0.459</b>	<b>0.163</b>	<b>0.105</b>	<b>0.513</b>

+5.4%

+5.7%

+0.7%

+3.6%

# Indexing Experts at the University of Pisa



URL: <https://wiser.d4science.org>




Search by Expertise




Search by Name

Search by Department

- ▷ ~1.5K Authors
- ▷ ~65K Documents (papers' abstracts)
- ▷ ~35K Research Topics
- ▷ More than 1K queries and ~2K profiles view in few months
- ▷ Currently used by UniPi's Technology Transfer Office

 Select the range of years to analyze.



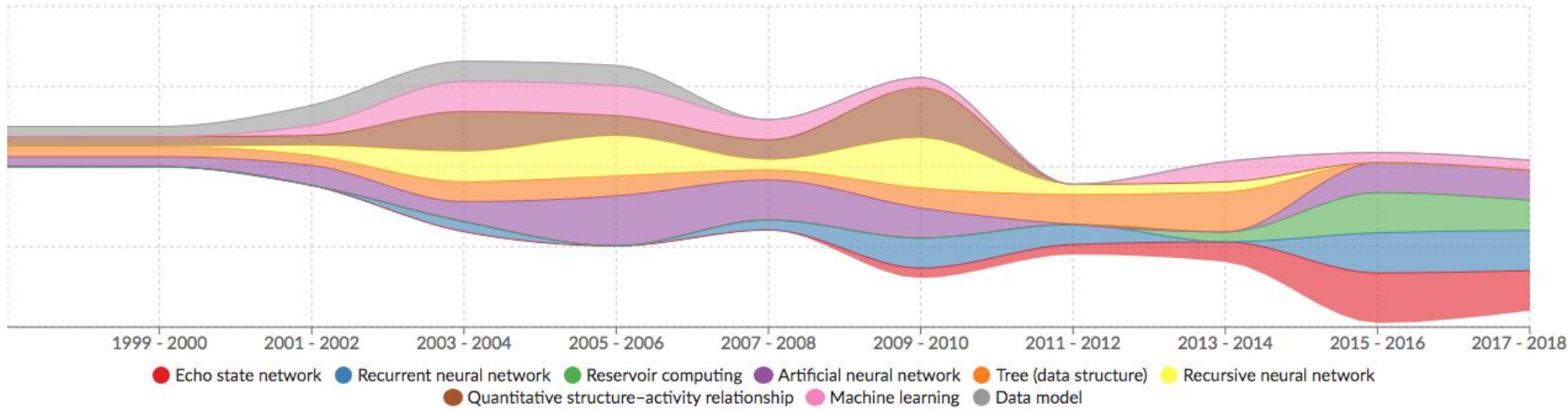
Ra...	Entity	Count	Doc. count	Years															Wiser score 
				1998	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	
1	<a href="#">Recurrent neural network</a>	26	15	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	
2	<a href="#">Artificial neural network</a>	44	24	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
3	<a href="#">Tree (data structure)</a>	24	17	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
4	<a href="#">Machine learning</a>	16	14	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
5	<a href="#">Recursive neural network</a>	16	16	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
6	<a href="#">Quantitative structure-activi...</a>	30	16	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
7	<a href="#">Echo state network</a>	14	13	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>		
8	<a href="#">Mathematical model</a>	39	28	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
9	<a href="#">Prediction</a>	19	17	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
10	<a href="#">Generative model</a>	9	7	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
11	<a href="#">Group representation</a>	22	14	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
12	<a href="#">Empiricism</a>	15	15	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
13	<a href="#">Data model</a>	8	8	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
14	<a href="#">Data set</a>	20	14	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		
15	<a href="#">Reservoir computing</a>	8	8	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		



Main Topics Main Areas **Stream Graph** Tag Cloud Publications Survey

Select range of years to examine.

Select how many topics for each range



# 4

## Future Research Directions

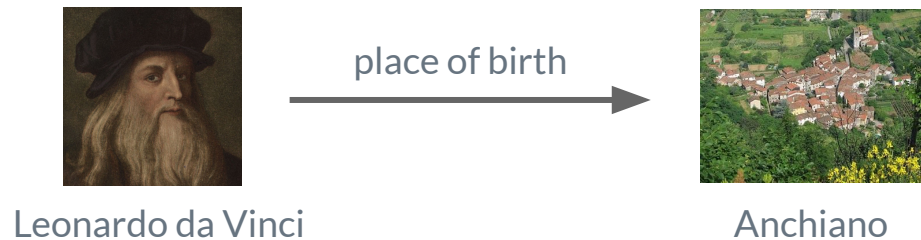


# Future Research Directions

## ▷ Entity Relatedness



- Apply our Two-Stage Framework over other KGs
- Extend it to *labels* associated to entities' *relationships*



# Conclusion and Future Directions

## ▷ Entity and Fact Salience

- Improve quality of Entity Salience annotations for NYT
- Entity Linking research will start focusing on efficiency (*best solutions are currently extremely slow, especially when applied over large-scale!*)
- Number of Applications for both Entity and Fact Salience
  - News Credibility (Popat, EMNLP'18)
  - KGs Construction (Nguyen, VLDB'18)
  - Facts Contextualization (Voskarides, SIGIR'18)

# Conclusion and Future Directions

## ▷ Expert Finding

- **Fine-grain Clustering** of Wikipedia entities for the visualization of groups of topics of an expert

Algorithms

Computer Science

Gzip

Burrows-Wheeler Transform

- Classical Clustering Algorithms generate one single cluster
  - Gzip and Burrows-Wheeler should actually belong to a domain-specific cluster
- Apply our graph-based profiling technique to **other domains**, e.g., recommendation systems, conversational AI

Thanks!  
**Any questions?**

# Contributions

1

## Entity Relatedness

A Two-Stage Framework for Computing  
Entity Relatedness in Wikipedia  
Marco Ponza, Paolo Ferragina and Soumen Chakrabarti



2

## Entity Salience

Document Aboutness via Sophisticated  
Syntactic and Semantic Features  
Marco Ponza, Paolo Ferragina and Francesco Piccinno



SWAT: A System for Detecting  
Salient Wikipedia Entities in Texts  
Marco Ponza, Paolo Ferragina and Francesco Piccinno



3

## Fact Salience

Facts That Matter  
Marco Ponza, Luciano Del Corro and Gerhard Weikum



4

## Expert Finding

WISER: A Semantic Approach for Expert Finding  
in Academia based on Entity Linking  
Paolo Cifariello, Paolo Ferragina and Marco Ponza



Information Systems 2019