

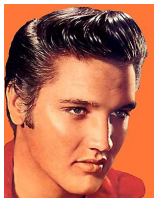
# Making knowledge bases more complete

Wikimedia Showcase 2019-12-18

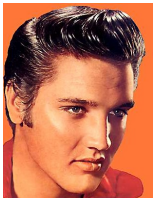
Fabian M. Suchanek

Télécom Paris University, France

# I am an Elvis Fan!

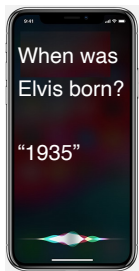


# We visited Elvis in New Zealand



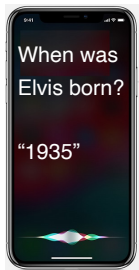
Elvis' lodge in Murchison

# New Technology: Knowledge Bases



Apple Siri

# New Technology: Knowledge Bases



Apple Siri

What is the capital  
of New Zealand?  
"Wellington"



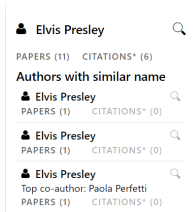
Amazon Echo

Discovered 6 kinase  
proteins that relate  
to cancer

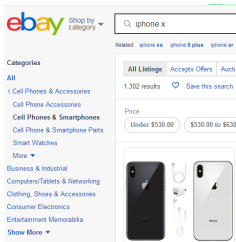


IBM Watson

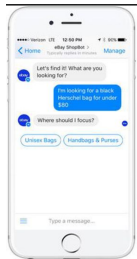
# New Technology: Knowledge Bases



Microsoft  
Academic



Ebay Knowledge Graph





Ebay Shopbot on  
Facebook  
Messenger

# New Technology: Knowledge Bases



Elvis Presley






**Elvis Presley** 

American singer

---

Available on

-  YouTube
-  Spotify
-  Deezer

▼ More music services

Elvis Aaron Presley, also known mononymously as Elvis, was an American singer, musician, and actor. Regarded as one of the most significant cultural icons of the 20th century, he is often referred to as the "King of Rock and Roll" or simply "the King". [Wikipedia](#)

**Born:** January 8, 1935, [Tupelo, Mississippi, United States](#)

**Died:** August 16, 1977, [Graceland, Memphis, Tennessee, United States](#)

Google Knowledge Graph

# New Technology: Knowledge Bases



Elvis Presley



Elvis Presley

American singer

Available on

YouTube

Spotify

Deezer

▼ More music services

Elvis Aaron Presley, also known mononymously as Elvis, was an American singer, musician, and actor. Regarded as one of the most significant cultural icons of the 20th century, he is often referred to as the "King of Rock and Roll" or simply "the King". [Wikipedia](#)

**Born:** January 8, 1935, Tupelo, Mississippi, United States

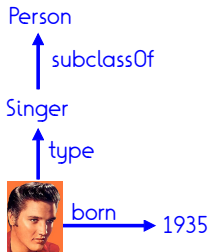
**Died:** August 16, 1977, Graceland, Memphis, Tennessee, United States

???

Google Knowledge Graph



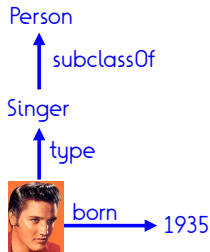
# Knowledge Bases



For us, a knowledge base (KB) is a graph, where the nodes are entities and the edges are relations.

(We do not distinguish T-Box and A-Box.)

# This talk



1. Constructing Knowledge Bases
2. Completing Knowledge Bases

# Extracting from Wikipedia

Elvis Presley



WIKIPEDIA  
The Free Encyclopedia

Elvis Presley was one of the best blah blah blub blah don't read this, listen to the speaker! blah blah blah blubl blah you are still reading this! blah blah blah blah blabbel blah



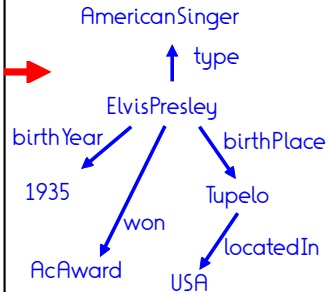
Born: 1935

In: Tupelo

...

Categories:

Rock&Roll, American Singers,  
Academy Award winners...





# YAGO: a large knowledge base



<http://yago-knowledge.org>  
open code and open data

Wikipedia + WordNet  
time and space  
10 languages  
100 relations  
100m facts  
10m entities  
95% accuracy  
used by DBpedia  
and IBM Watson

New version  
in preparation!

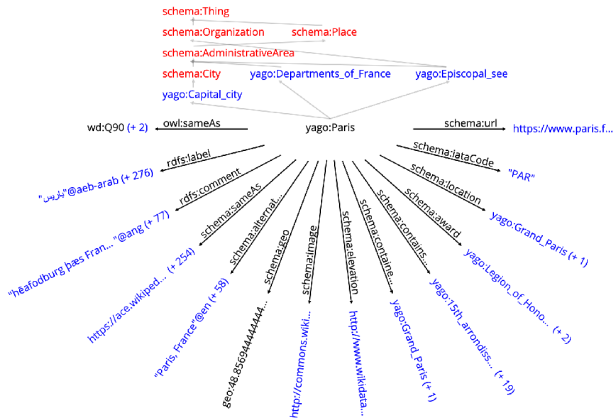


[WWW'07, JWS'08, WWW'11 demo, AIJ'13, WWW'13 demo, CIDR'15, ISWC'16]

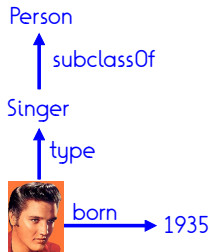


# YAGO 4: A "reason-able" KB

YAGO 4 combines schema.org + Wikidata + constraints

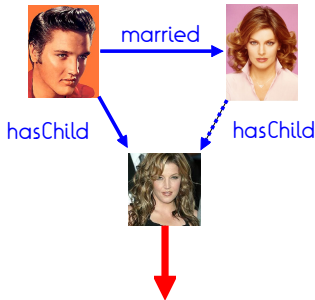


# This talk



1. Constructing Knowledge Bases ✓
2. Completing Knowledge Bases

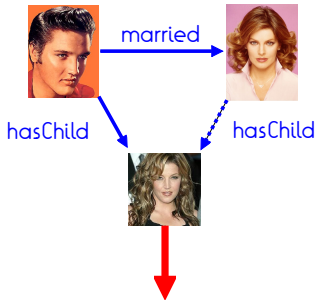
# Incompleteness: Concrete facts



$$\text{married}(x,y) \wedge \text{hasChild}(x,z) \Rightarrow \text{hasChild}(y,z)$$



# Incompleteness: Concrete facts



$$\text{married}(x,y) \wedge \text{hasChild}(x,z) \Rightarrow \text{hasChild}(y,z)$$

But: Rule mining needs counter examples  
and RDF ontologies are positive only

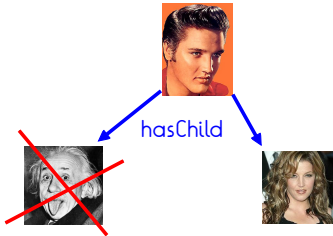
# Partial Completeness Assumption



Assumption:

If we know  $r(x, y_1), \dots, r(x, y_n)$ ,  
then all other  $r(x, z)$  are false.

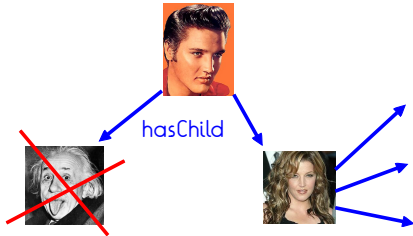
# Partial Completeness Assumption



Assumption:

If we know  $r(x, y_1), \dots, r(x, y_n)$ ,  
then all other  $r(x, z)$  are false.

# Partial Completeness Assumption



Assumption:

If we know  $r(x, y_1), \dots, r(x, y_n)$ ,  
then all other  $r(x, z)$  are false.

# AMIE finds rules in knowledge bases



AMIE

(2min)

$$r(x, y) \wedge r'(z, y) \Rightarrow r''(x, z)$$

AMIE is based on an efficient in-memory database implementation.

Caveat: rules cannot  
predict the unknown  
with high precision

# AMIE finds rules in knowledge bases



AMIE



$type(x, pope) \Rightarrow$   
 $diedIn(x, Rome)$



[WWW 2013, VLDB journal 2015]

New version  
in preparation!



# Incompleteness: Existence of facts



marriedTo



the quality of YAGO w  
ls a precision of 95%, as  
uks to our brilliant algori

# Incompleteness: Existence of facts



marriedTo

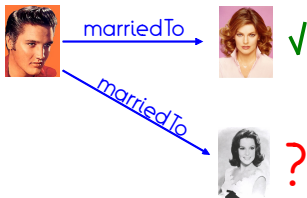


the quality of YAGO w  
is a precision of 95%, as  
uks to our brilliant algori





# Incompleteness: Existence of facts



the quality of YAGO w  
ls a precision of 95%, as  
iks to our brilliant algori



Given a subject  $s$  and  
a relation  $r$ , do we know  
all  $o$  with  $r(s, o)$  ?

# Signals for Incompleteness



Closed World Assumption

Partial Completeness Assumption

Popularity oracle

No-change oracle

Star-pattern oracle

Class-oracle

AMIE oracle: Learn rules such as

$moreThan_1(x, hasParent) \Rightarrow complete(x, hasParent)$



[>details](#)

# Signals for Incompleteness (F1)

Relation	CWA	PCA	card <sub>2</sub>	Popularity	No change	Star	Class	AMIE
diedIn	60%	22%	—	4%	15%	50%	<b>99%</b>	96%
directed	40%	96%	19%	7%	71%	0%	0%	<b>100%</b>
graduatedFrom	89%	4%	2%	2%	10%	89%	<b>92%</b>	87%
hasChild	71%	1%	1%	2%	13%	40%	<b>78%</b>	<b>78%</b>
hasGender	78%	<b>100%</b>	—	2%	—	86%	95%	<b>100%</b>
hasParent*	1%	54%	<b>100%</b>	—	—	0%	0%	<b>100%</b>
isCitizenOf*	4%	98%	11%	1%	4%	10%	5%	<b>100%</b>
isConnectedTo	87%	34%	19%	—	—	68%	88%	<b>89%</b>
isMarriedTo*	55%	7%	0%	3%	12%	37%	<b>57%</b>	46%
wasBornIn	28%	<b>100%</b>	—	5%	8%	0%	0%	<b>100%</b>



Relation	CWA	PCA	card <sub>2</sub>	Popularity	Star	Class	AMIE
alma_mater	<b>90%</b>	14%	5%	1%	87%	87%	87%
brother	93%	1%	—	1%	94%	<b>96%</b>	<b>96%</b>
child	70%	1%	—	1%	<b>79%</b>	72%	73%
country_of_citizenship*	42%	97%	10%	3%	0%	0%	<b>98%</b>
director	81%	<b>100%</b>	—	3%	94%	89%	<b>100%</b>
father*	5%	<b>100%</b>	6%	9%	89%	8%	<b>100%</b>
mother*	3%	<b>100%</b>	3%	10%	67%*	5%	<b>100%</b>
place_of_birth	53%	<b>100%</b>	7%	5%	55%	0%	<b>100%</b>
place_of_death	89%	35%	1%	2%	81%	81%	<b>96%</b>
sex_or_gender	81%	<b>100%</b>	6%	3%	92%	91%	<b>100%</b>
spouse*	<b>57%</b>	7%	—	1%	54%	54%	55%



\* = biased training sample

# Incompleteness: Existence of facts



marriedTo



marriedTo



AMIE can predict incompleteness

- bornIn: 100% F1-measure
- diedIn: 96%
- directed: 100%
- graduatedFrom: 87%
- hasChild: 78%
- isMarriedTo: 46%
- ... and more.



[WSDM 2017]

>rep&married

>married

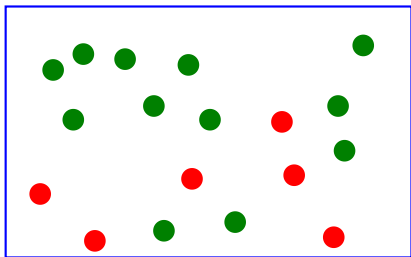
# Are all people married?



# Are all people married?



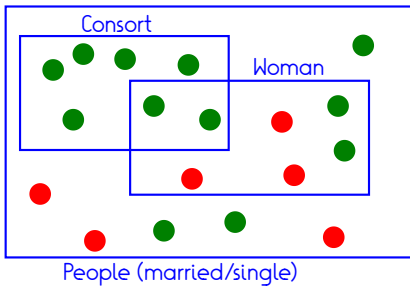
# Are all people married?



Real World

People (married/single)

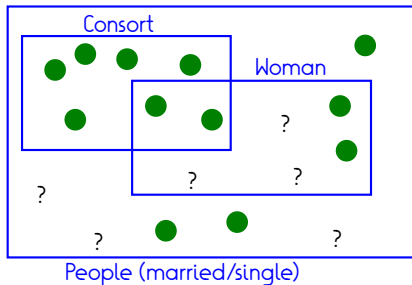
# Are all people married?



Real World

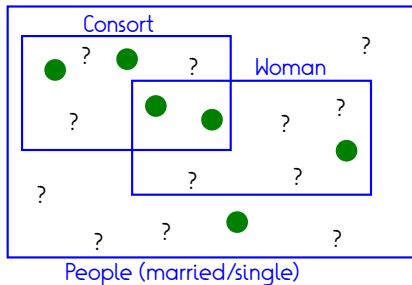


# Are all people married?



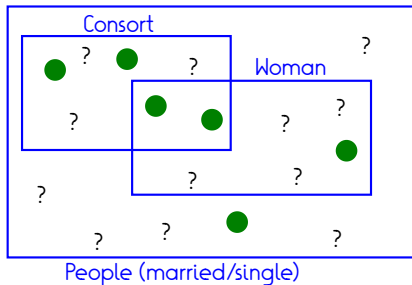
Knowledge base  
under the  
Open World Assumption

# Are all people married?



Knowledge base  
under the  
Open World Assumption  
and incompleteness

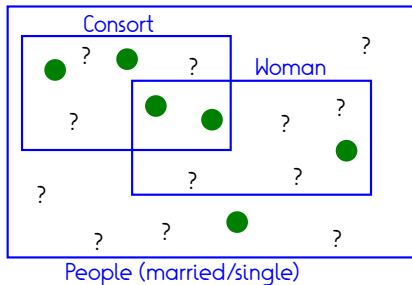
# Are all people married?



Knowledge base  
under the  
Open World Assumption  
and incompleteness

Baseline 1: Obligatory if all instances have it

# Are all people married?

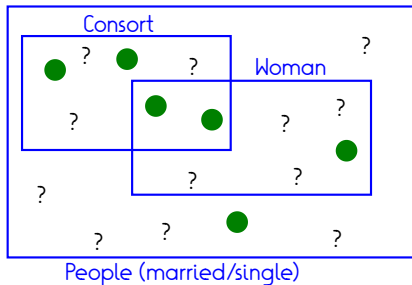


Knowledge base  
under the  
Open World Assumption  
and incompleteness

Baseline 1: Obligatory if all instances have it ✗

Baseline 2: Obligatory if at least n% of instances have it

# Are all people married?



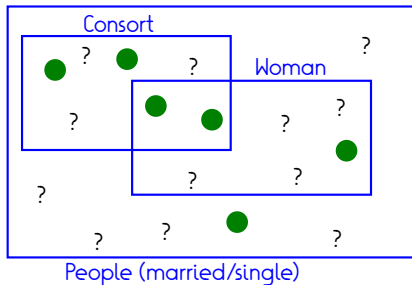
Knowledge base  
under the  
Open World Assumption  
and incompleteness

Baseline 1: Obligatory if all instances have it ✗

Baseline 2: Obligatory if at least n% of instances have it => Woman ✗

Baseline 3: Obligatory if all instances that have it fall in the class

# Are all people married?



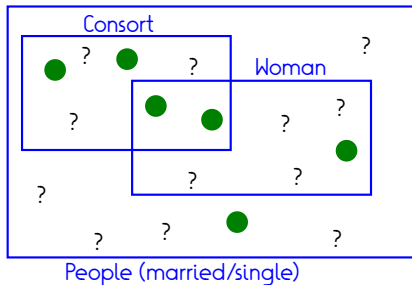
Knowledge base  
under the  
Open World Assumption  
and incompleteness

Baseline 1: Obligatory if all instances have it ✗

Baseline 2: Obligatory if at least n% of instances have it => Woman ✗

Baseline 3: Obligatory if all instances that have it fall in the class  
=> Person ✗

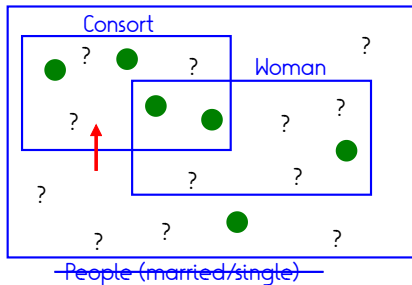
# Are all people married?



Knowledge base  
under the  
Open World Assumption  
and incompleteness

Theorem: If the KB is sampled randomly uniformly from the real world, and if the density of an attribute changes when we go into an intersecting class, then the attribute cannot be obligatory.

# Are all people married?

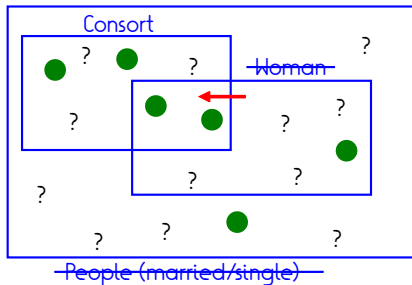


Knowledge base  
under the  
Open World Assumption  
and incompleteness

Theorem: If the KB is sampled randomly uniformly from the real world, and if the density of an attribute changes when we go into an intersecting class, then the attribute cannot be obligatory.



# Are all people married?

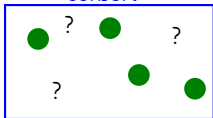


Knowledge base  
under the  
Open World Assumption  
and incompleteness

Theorem: If the KB is sampled randomly uniformly from the real world, and if the density of an attribute changes when we go into an intersecting class, then the attribute cannot be obligatory.

# Determining obligatory attributes

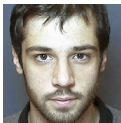
Consort



We can predict obligatory attributes of classes with up to 80% precision (at 40% recall).

Caveat:

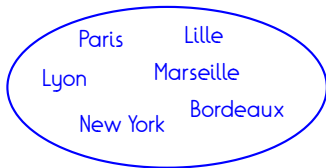
We do not actually predict, but exclude.



[WWW 2018]

# Incompleteness: Missing entities

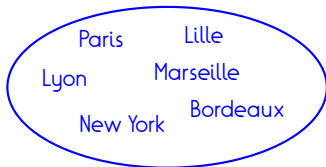
We have the following cities in our knowledge base:



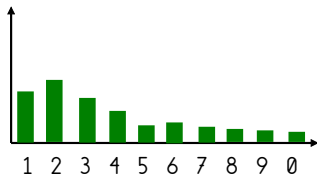
Are there any cities missing?

# Incompleteness: Missing entities

We have the following cities in our knowledge base:



Are there any cities missing?



- 1) Take the number of inhabitants of each city
- 2) Take the first digit
- 3) Plot the number of cities per first digit

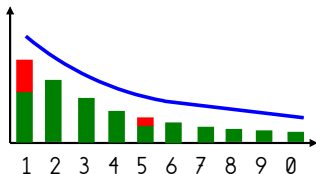
# Incompleteness: Missing entities

Benford's law says that the first digit  $d$  appears with probability

$$\log_{10}\left(1 + \frac{1}{d}\right)$$

=> We can give a minimum numbers of cities that are missing to make the distribution representative of the real world.

(For other classes, we can learn a parameter for a variant of the law.)



[ISWC 2018]




# This talk



1. Constructing Knowledge Bases ✓
2. Completing Knowledge Bases ✓

# Is Elvis dead?



## Elvis Presley

American singer

Available on

- YouTube
- Spotify
- Deezer
- More music services

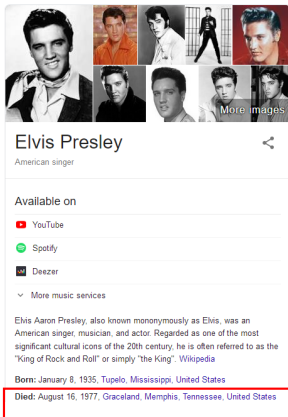
Elvis Aaron Presley, also known mononymously as Elvis, was an American singer, musician, and actor. Regarded as one of the most significant cultural icons of the 20th century, he is often referred to as the "King of Rock and Roll" or simply "the King". [Wikipedia](#)

Born: January 8, 1935, Tupelo, Mississippi, United States

Died: August 16, 1977, Graceland, Memphis, Tennessee, United States

???

# Is Elvis dead?



A profile card for Elvis Presley. At the top is a grid of eight images: a large black and white portrait of him smiling, and seven smaller images showing him in various styles (orange shirt, white shirt, black jumpsuit, red shirt, black shirt, white shirt, and black shirt). Below the grid is the text 'Elvis Presley' and 'American singer'. Underneath is a section 'Available on' with icons for YouTube, Spotify, and Deezer, followed by a link 'More music services'. A paragraph of text describes him as 'Elvis Aaron Presley, also known mononymously as Elvis, was an American singer, musician, and actor. Regarded as one of the most significant cultural icons of the 20th century, he is often referred to as the "King of Rock and Roll" or simply "the King". Wikipedia'. At the bottom, it says 'Born: January 8, 1935, Tupelo, Mississippi, United States' and 'Died: August 16, 1977, Graceland, Memphis, Tennessee, United States'. The 'Died' line is highlighted with a red border.

Elvis Presley

American singer

Available on

YouTube

Spotify

Deezer

More music services

Elvis Aaron Presley, also known mononymously as Elvis, was an American singer, musician, and actor. Regarded as one of the most significant cultural icons of the 20th century, he is often referred to as the "King of Rock and Roll" or simply "the King". Wikipedia

Born: January 8, 1935, Tupelo, Mississippi, United States

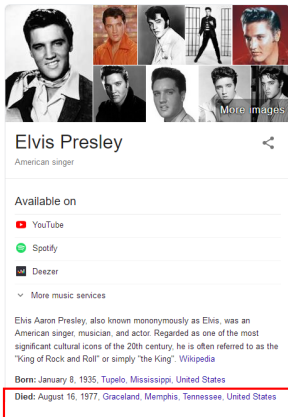
Died: August 16, 1977, Graceland, Memphis, Tennessee, United States



???



# Is Elvis dead?



A profile card for Elvis Presley. At the top is a grid of eight images: a large black and white portrait of him smiling, and seven smaller images showing him in various poses and outfits. Below the images is the text 'Elvis Presley' and 'American singer'. Underneath is a section 'Available on' with icons for YouTube, Spotify, and Deezer, and a link to 'More music services'. A paragraph of text describes him as an American singer, musician, and actor, known as 'the King of Rock and Roll'. At the bottom, it lists his birth and death information. The death information is highlighted with a red box.

**Elvis Presley**  
American singer

Available on

- YouTube
- Spotify
- Deezer
- More music services

Elvis Aaron Presley, also known mononymously as Elvis, was an American singer, musician, and actor. Regarded as one of the most significant cultural icons of the 20th century, he is often referred to as the "King of Rock and Roll" or simply "the King". [Wikipedia](#)

Born: January 8, 1935, Tupelo, Mississippi, United States

**Died: August 16, 1977, Graceland, Memphis, Tennessee, United States**



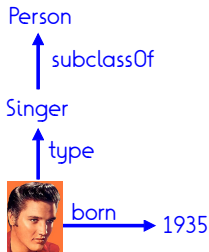
100m statements

95% accuracy

-> 5m wrong statements

???

# Knowledge Bases



1. Constructing Knowledge Bases ✓
2. Completing Knowledge Bases ✓