



# Towards a Vecsigrafo Portable Semantics in Knowledge-based Text Analytics

Ronald Denaux & José Manuel Gómez Pérez

HSSUES– Oct. 21st, 2017

# The Cognitive Chasm

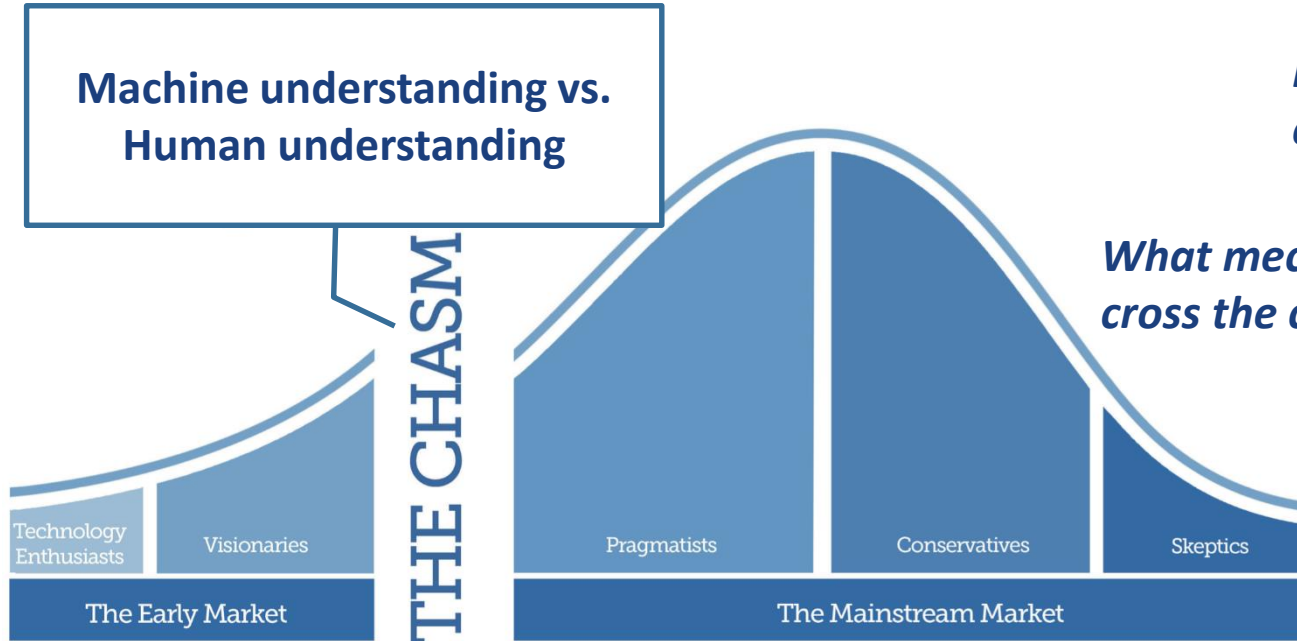
*How can humans and AI interact with and understand each other?*

*Is this possible or are they cognitively disconnected?*

**Machine understanding vs.  
Human understanding**

*What mechanisms are needed to cross the cognitive chasm?*

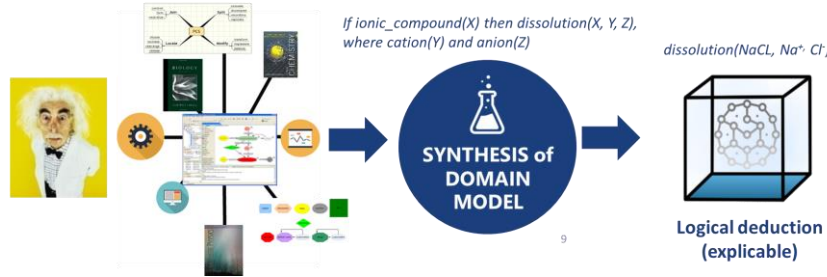
*How can knowledge representation be both flexible, scalable, deep and logical?*



# Pros and cons of structured knowledge

## PROS

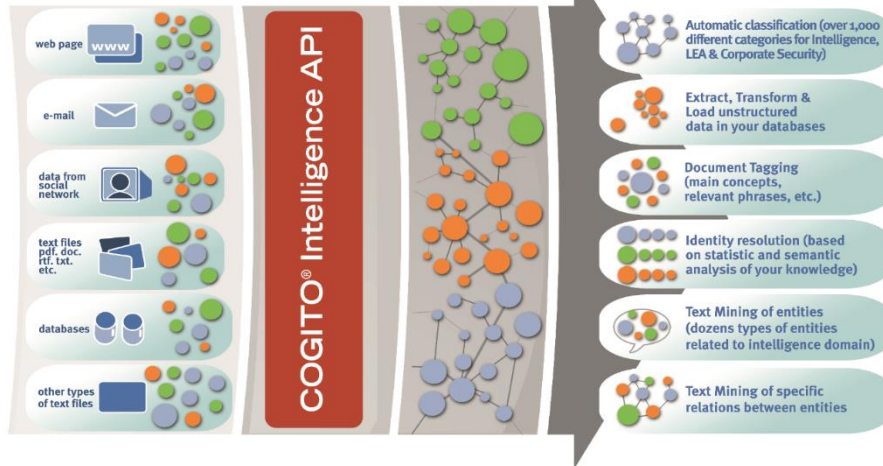
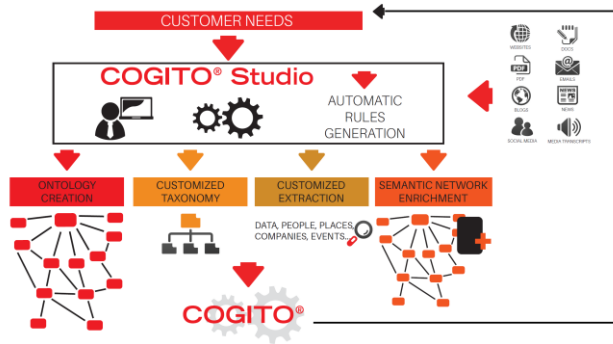
- Humans have a rich understanding of the domain, resulting in detailed, expressive models
- Underlying formalisms support logical explanations
- Reasonable response times
- Tooling can optimize cost, enabling user-entered knowledge



## CONS

- Requires a considerable amount of well trained, centralized labor to manually encode knowledge
- Lacks scalability with large corpora and still costly due to humans in the loop
- Possible bias, hard to generalize
- Brittleness

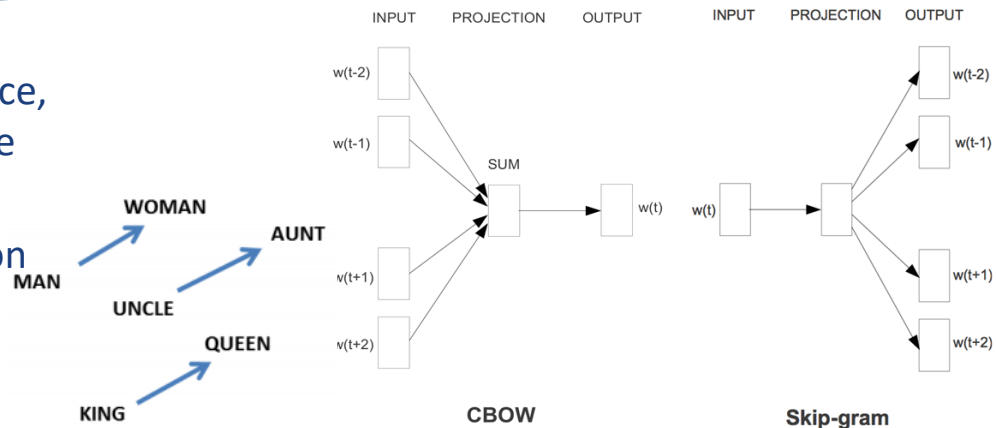
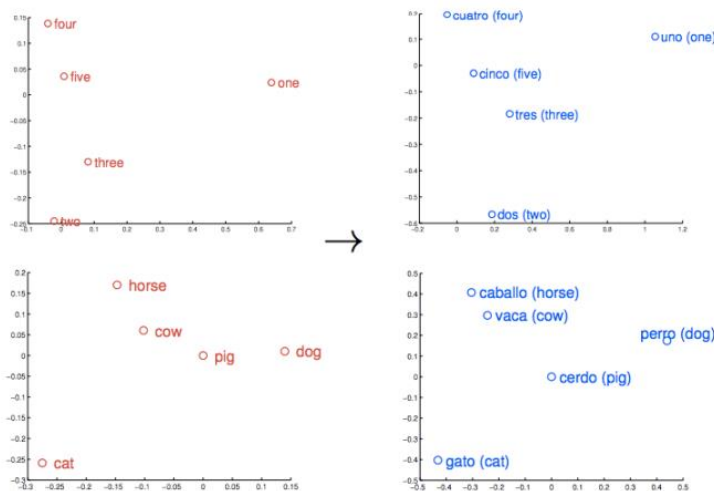
# Structured knowledge (Sensigrafo)



- **Sensigrafo, a knowledge graph** containing word definitions, related concepts and linguistic information
- Main entities include **syncons** (concepts), **lemmas** (canonical representation of a word) and **relations** (properties, taxonomical, polysemy, synonymy...)
  - 301,582 syncons
  - 401,028 lemmas
  - 80+ relation types that yield ~2.8 million links
- **Internal representation** that leverages external resources, both general and domain-specific
- **Word-sense disambiguation**, based on the context of a word in Sensigrafo
- **Categorization and extraction** supported through **Sensigrafo plus lexical-syntactic rules**

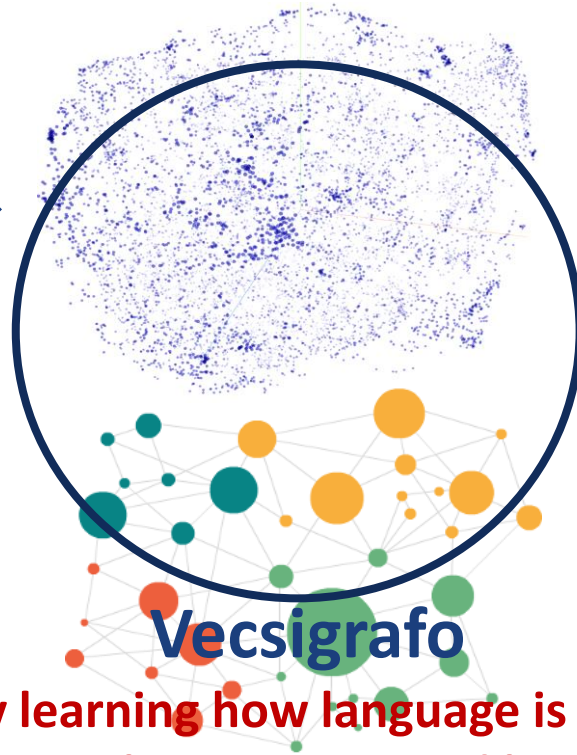
# Building multiple language models

- **Word2vec** represents words in a vector space, making natural language computer-readable
- Neural word embeddings enable **word similarity, analogy and relatedness** based on vector arithmetic (cosine similarity)
- **Essential property: Semantic portability**

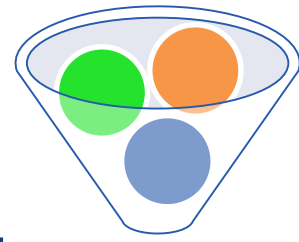


Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

# Towards Natural Learning at Expert System

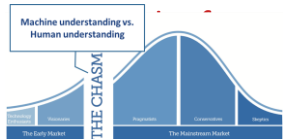


- Knowledge encoded in the mind of the expert
- Structured knowledge base
- Good for logical deduction and explanation
- Deep, but rigid and brittle
- Human is a bottleneck: hand-engineered features and powerful modeling tools needed



- Knowledge embedded in document corpora
- Broad, flexible, scalable
- Good for POS tagging, parsing, semantic relatedness
- Statistic induction, not logical explanation

**Automatically learning how language is used in real life and materializing that in structured knowledge graphs**

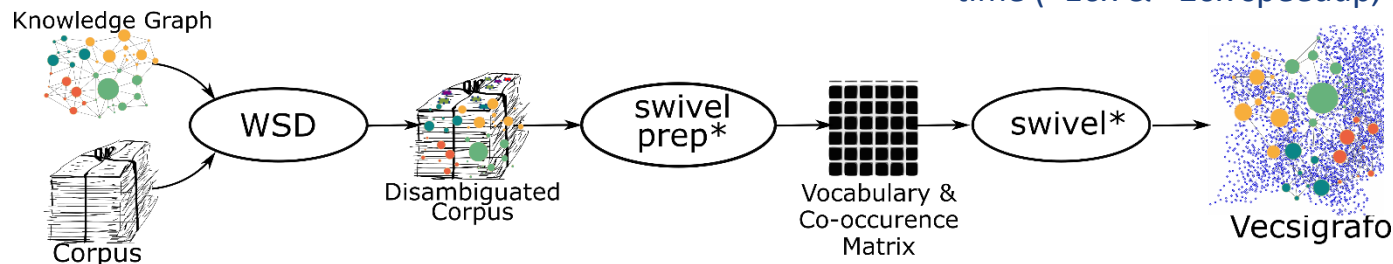


# Vecsigrafo – Putting it all together

Vocab elements	EN-grafo		ES-grafo	
	Sensi	Vecsi	Sensi	Vecsi
Lemmas	398	80	268	91
Concepts	300	67	226	52
<b>Total</b>	<b>698</b>	<b>147</b>	<b>474</b>	<b>143</b>

Corpus	Sentences	Spanish words	English words
Euparl	1,965,734	51,575,748	49,093,806
<b>UN.en-es</b>	<b>21,911,121</b>	<b>678,778,068</b>	<b>590,672,799</b>

- **Two parallel corpora**, focused on English and Spanish (**Europarl** and **UN**)
- **Meaning extracted from corpora and related to Sensigrafo** (**21%** and **30%** Sensigrafo covered, resp.)
- **Tokenized, lemmatized and disambiguated** with COGITO
- Learned **monolingual joint word-concept models** and a **(non-linear) transformation** between vector spaces for crosslinguality
- **Deeplearning4j** with Skip-gram, minFreq 10, vector dimensionality 400
- **TensorFlow** and **Swivel** for better vectorization time (~16x & ~20x speedup, 80 epochs)



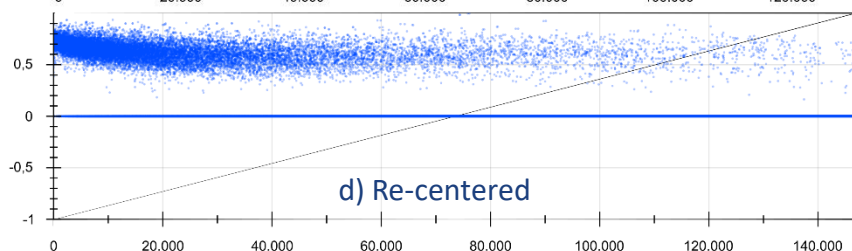
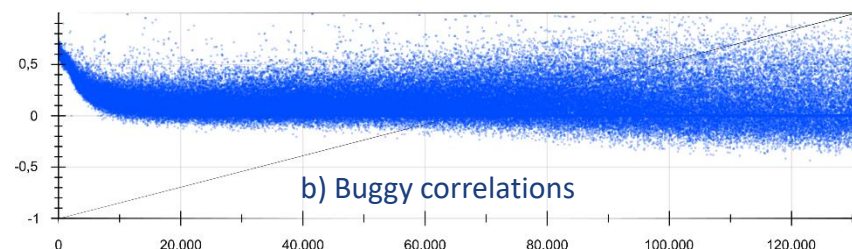
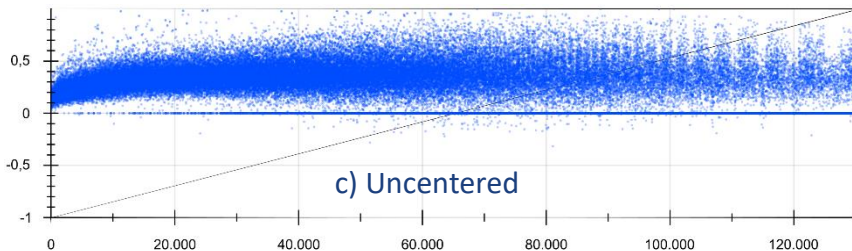
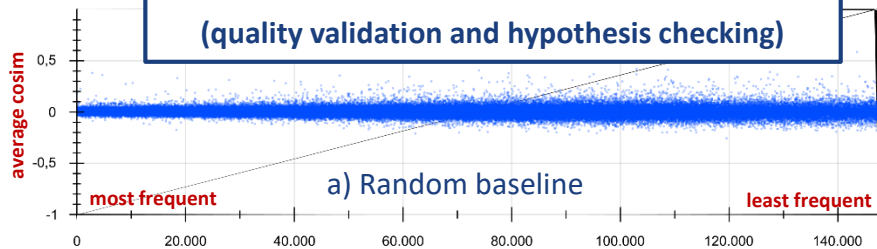


# Vecsigrafo - Evaluation

Model	WSim	WSrel	Simlex999	Rarewords	Simverb
SotA 2015	79.4	70.6	43.3	50.8	n/a
Swivel	74.8	61.6	40.3	48.3	62.8
Swivel <sub>UN, en</sub>	58.8	45.0	18.3	37.8	15.3
<b>Vecsigrafo<sub>UN, en</sub></b>	<b>47.6</b>	<b>24.1</b>	<b>12.4</b>	<b>30.8</b>	<b>13.2</b>

## Word Prediction Plots

(quality validation and hypothesis checking)



- Corpus size and distribution matters
- ~~Overall performance equivalent at lemma level (Swivel, same corpus)~~
- ~~Including concepts has a cost~~
- Visual inspection (t-SNE, PCA) and manual (relatedness, analogy...)
- Further insight needed

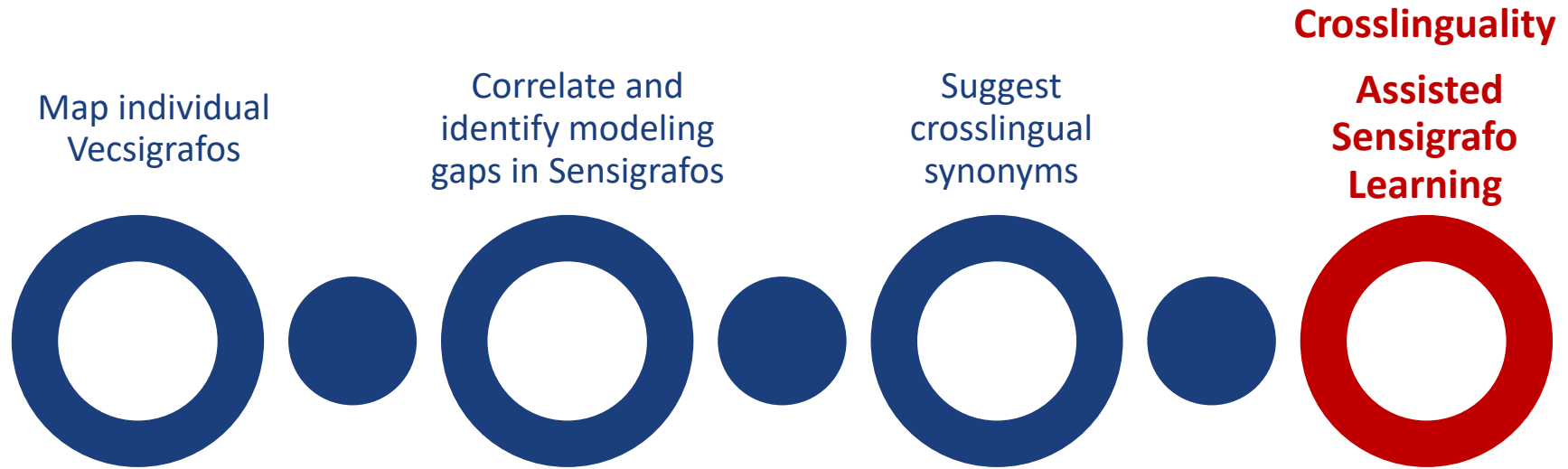


# Vecsigrafo – Word Similarity Redux

Model	WSim	WSrel	Simlex999	Rarewords	Simverb
SotA 2015	79.4	70.6	43.3	50.8	62.8
Swivel	74.8	61.6	40.3	48.3	n/a
<b>Swivel<sub>UN, en</sub></b>	<b>58.8</b>	45.0	18.3	37.8	15.3
<b>Swivel<sub>UN, en</sub> recentered</b>	57.7	<b>47.2</b>	<b>21.3</b>	<b>39.2</b>	<b>17.0</b>
<del>Vecsigrafo<sub>UN, en</sub></del>	<del>47.6</del>	<del>24.1</del>	<del>12.4</del>	<del>30.8</del>	<del>13.2</del>
<b>Vecsigrafo<sub>UN, en</sub></b>	<b>69.9</b>	<b>51.6</b>	38.2	<b>50.3*</b>	<b>30.6</b>
<b>Vecsigrafo<sub>UN, en</sub> recentered</b>	59.3	43.0	<b>42.4</b>	49.3	30.4
<b>Vecsigrafo<sub>UN, en</sub> NN aligned to es</b>	65.8	45.3	39.2	49.3	28.5

- Better than swivel for same corpus
- Effect of **recentering**
- Effect of **aligning to Spanish**
- **Further insight needed**
  - How similar are two vecsigrafos?
  - Which relations are inferred?
  - How are relations encoded in the embedding space?

# Vecsigrafo – Application Roadmap



Fast internationalization at Expert System (EU, US, LATAM) and growing customer needs in 14 languages

# Mapping and correlation

- **Mapping vector spaces in different languages:** Linear transformation suggested by (Mikolov, 2013) produced poor results. **Non-linear** transformation using NNs: **hit@5 = 0.78 and 90% semantic relatedness**
- Manual inspection showed **only 28% exact correspondence EN→ES**, due to volume (75K concepts less in Spanish Sensigrafo) and strategic modeling decisions
- **How to address the gap?**

Alignment performance

Method	Nodes	hit@5
TM	n/a	0.36
NN2	4K	0.61
NN2	5K	0.68
<b>NN2</b>	<b>10K</b>	<b>0.78</b>
NN3	5K	0.72

Manual inspection EN→ES

	in dict.	out dict.
#concepts	46	64
<b>hit@5</b>	<b>0.72</b>	<b>0.28</b>
no concept <sub>ES</sub>	2	33

# Examples

## “Financing” (EN→ES)

@	lemma/syncon	cosim	comment
1	financing	0.96	lemma
2	finance	0.85	lemma
3	funding	0.80	lemma
4	en#178501: adverb for financing	0.79	
5	en#75764: verb for fund, finance	0.76	
6	es#126922: noun financiación,	0.75	synonym



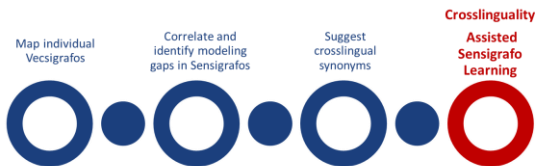
## “Scrap value” (EN→ES)

@	lemma/syncon	cosim	comment
1	salvage value	0.92	lemma
2	scrap value	0.90	lemma
3	replacement cost	0.72	
4	en#57338: replacement cost	0.72	
30	en#195309: reduced price, sale	0.61	
??	precio de compra	0.48	
??	es#20836: cambio, valor comercial	0.23	



## “PYME” (ES→EN)

@	lemma/syncon	cosim	comment
1	es#92662: pequeña y mediana empresa PYME	0.99	
2	en#2739337:SME	0.81	synonym
3	sme	0.81	synonym
4	mediana empresa	0.79	narrower
5	es#307734: mediana empresa	0.78	narrower

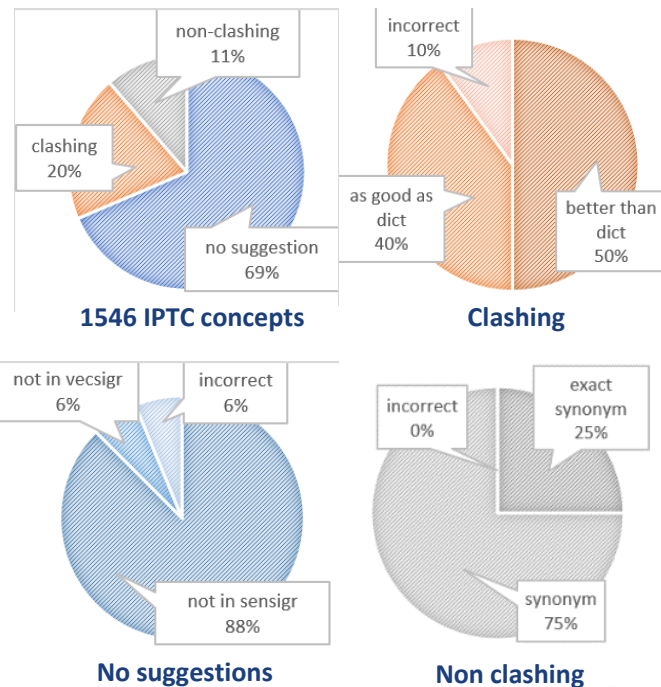


# Crosslingual synonym suggester

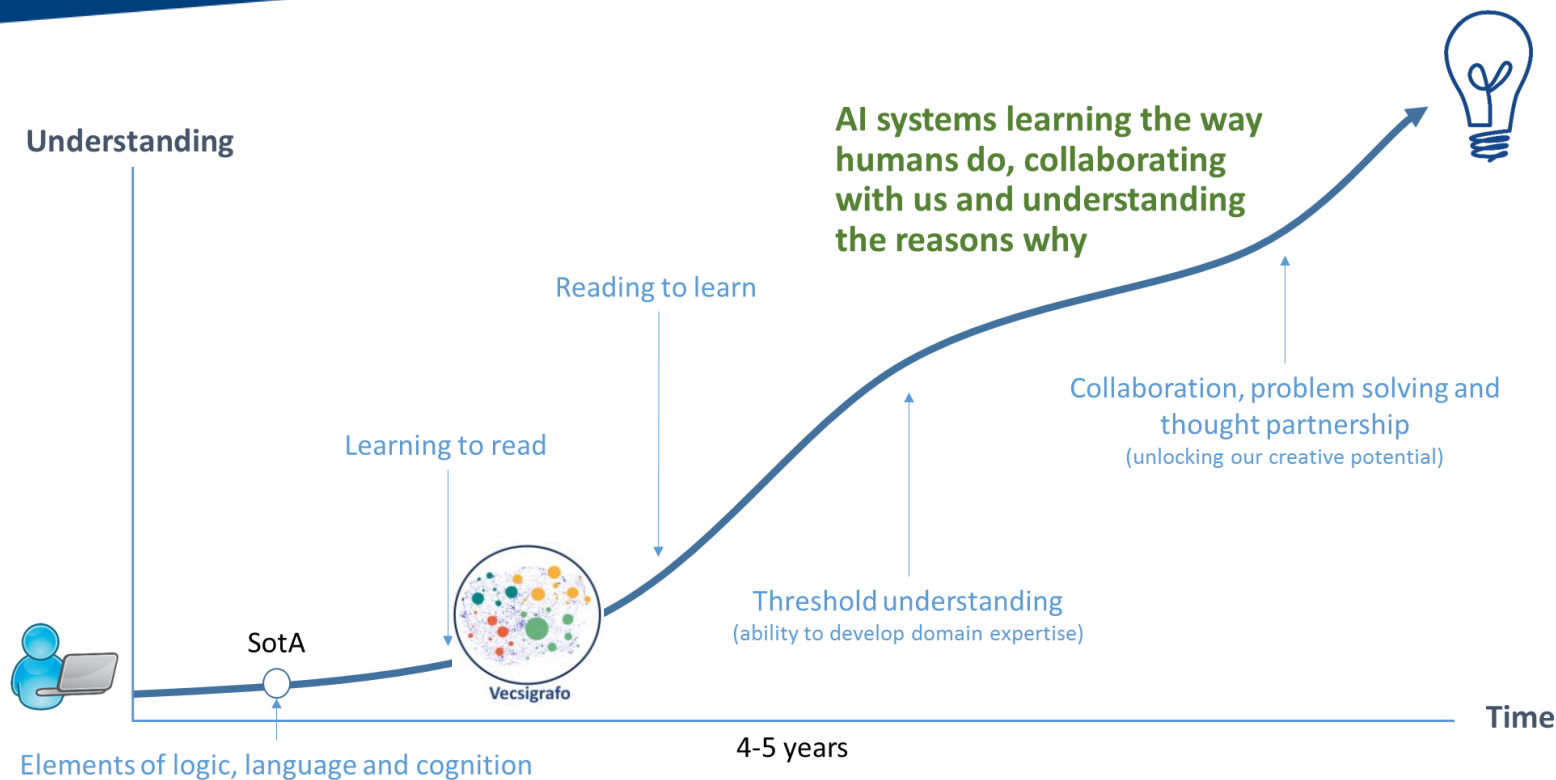
## Manual inspection EN→ES (1546 concepts, IPTC)

Combines features from bilingual vecsgrafo, the target and source Sensigrafos and a dictionary (PanLex)

1. For each concept in the source language, **find the  $n$  nearest concepts** in the target language that match grammar type (noun, verb, adjective, etc.)
2. For each candidate, **calculate hybrid features** (lemma translation, glossa similarity, cosine similarity, shared hypernyms and domains)
3. **Combine into a single score and rank**
4. Check if suggested synonym candidate is **already mapped** to a different concept and compare
5. Suggestion made if score is over a threshold



# Wrapping up



From David Ferrucci, 2016

Ronald Denaux  
Senior Researcher  
rdenaux@expertsystem.com

Jose Manuel Gomez-Perez  
Director R&D  
jmgomez@expertsystem.com



*Denaux R, Gomez-Perez JM. **Towards a Vecsigrafo: Portable Semantics in Knowledge-based Text Analytics.** To appear in proceedings of the Intl. Workshop on Hybrid Statistical Semantic Understanding and Emerging Semantics (HSSUES), collocated with the 16<sup>th</sup> Intl. Semantic Web Conference (ISWC), Vienna, 2017.*



[linkedin.com/company/expert-system](https://www.linkedin.com/company/expert-system)



[twitter.com/Expert\\_System](https://twitter.com/Expert_System)



[info@expertsystem.com](mailto:info@expertsystem.com)





# Correlation calculation

Develop an indicative list of advisory and conciliatory measures to encourage full compliance;



**Tokenize & WSD**

en#67083|develop en#89749|indicative en#113271|list en#88602|advisory en#85521|conciliatory en#33443|measure en#77189|encourage en#84127|full en#4941|compliance



**Correlation for en\_lem\_list (window 2, harmonic weight)**

token	Distance	weight
en#67083	2	$\frac{1}{2}$
develop	2	$\frac{1}{2}$
en#89749	1	1
indicative	1	1
en#113271	0	1

token	Distance	weight
list	0	1
en#88602	1	1
advisory	1	1
en#85521	2	$\frac{1}{2}$
conciliatory	2	$\frac{1}{2}$