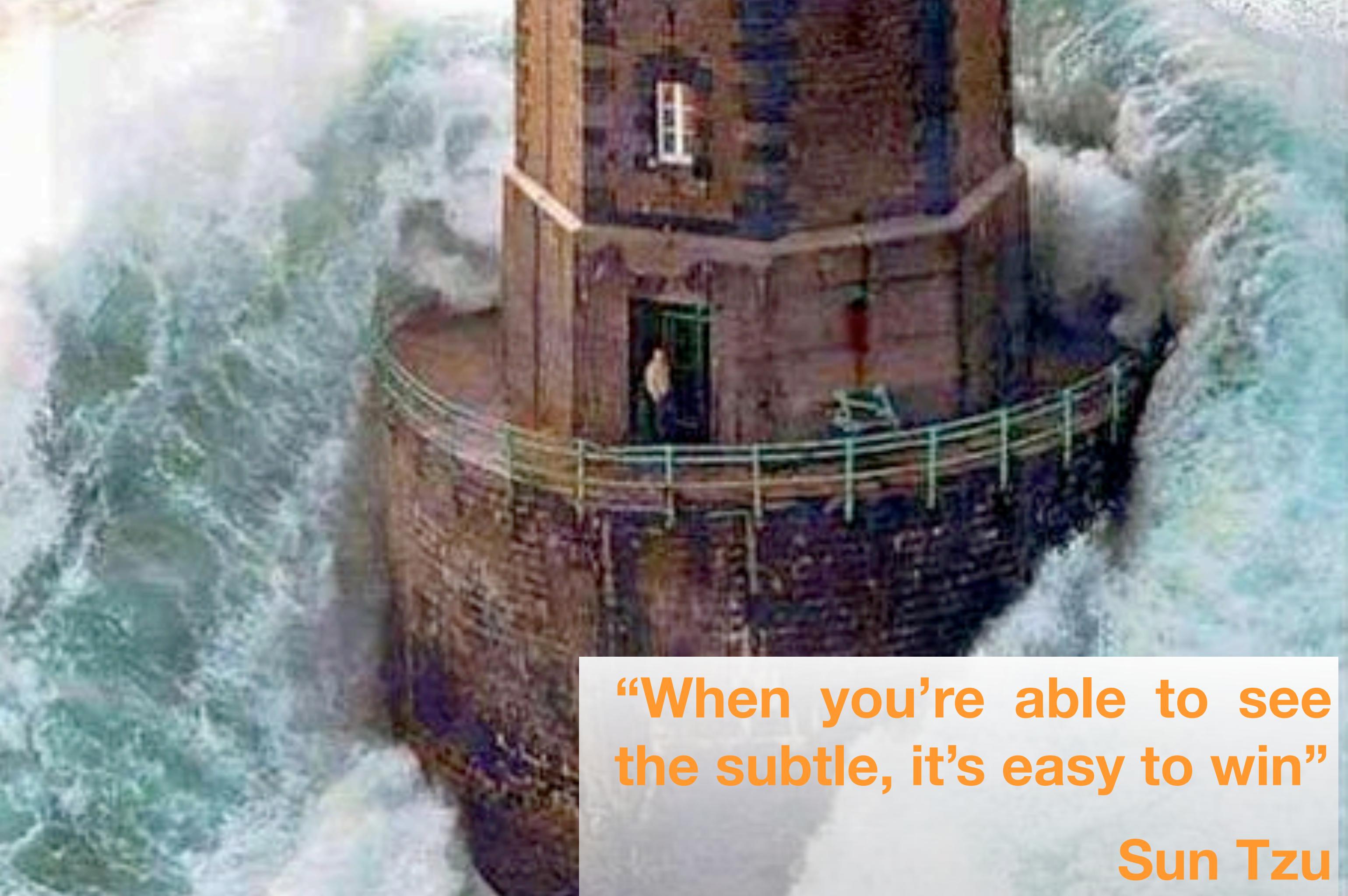


AUTORITAS



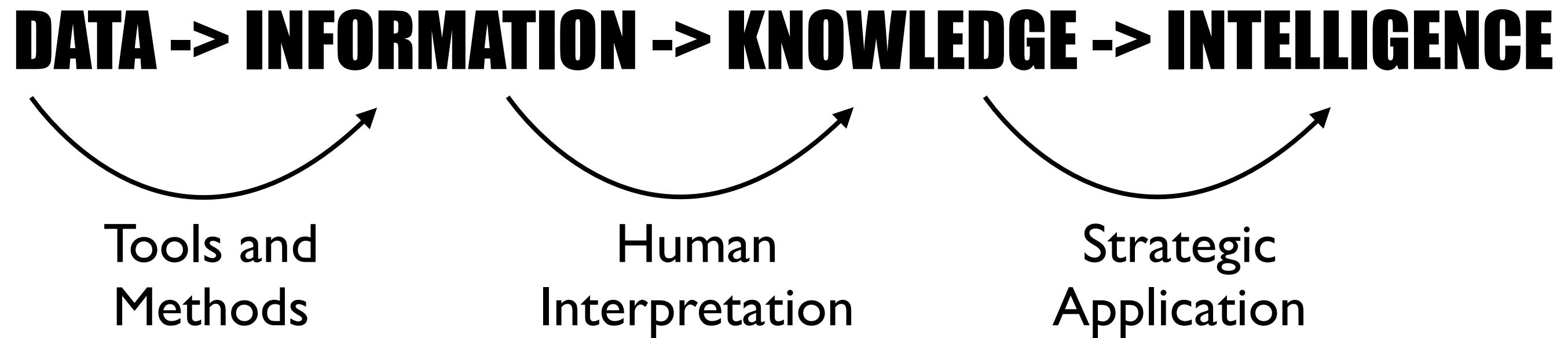
Smart Listening



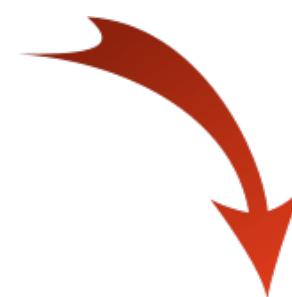
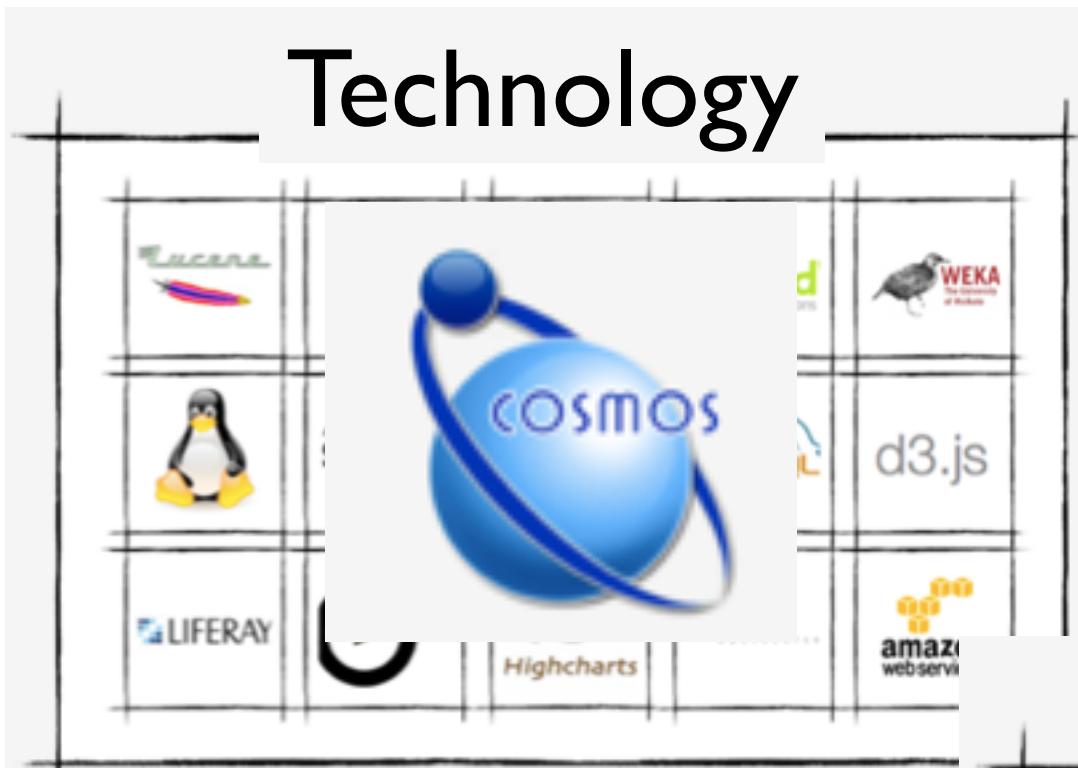
**“When you’re able to see
the subtle, it’s easy to win”**

Sun Tzu

Smart Listening

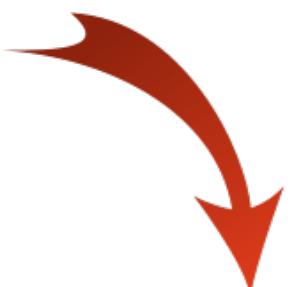
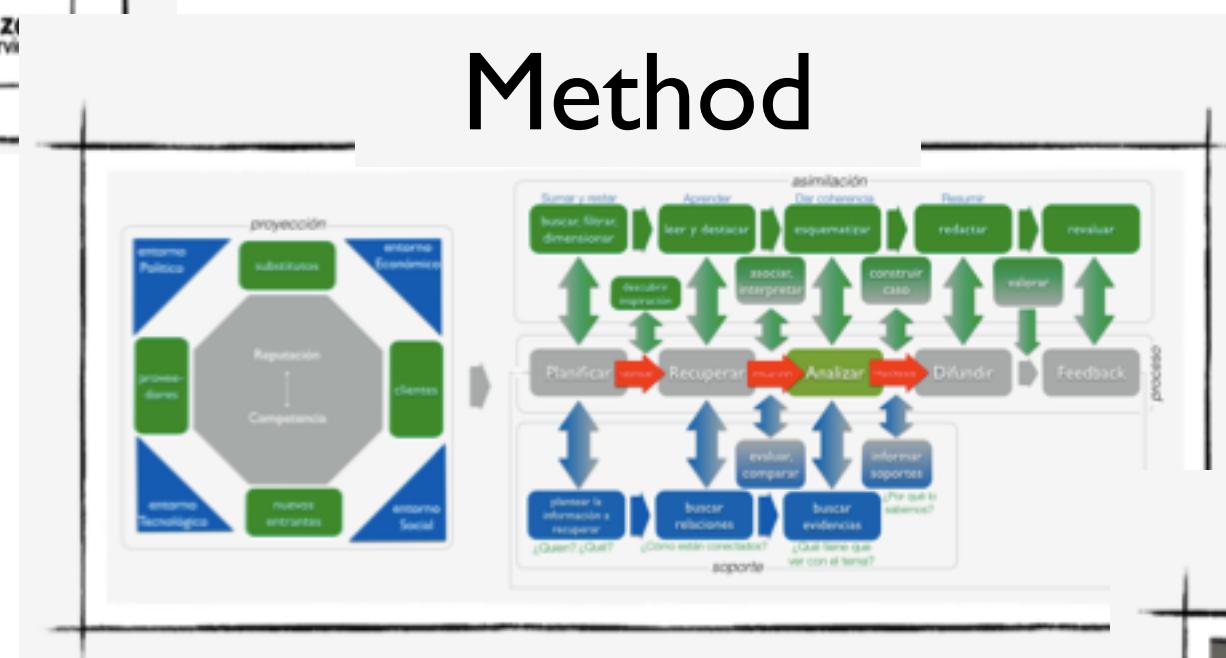


Technology

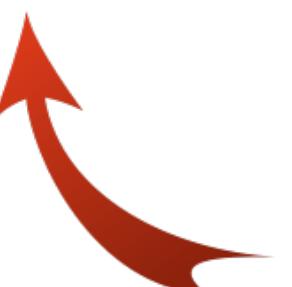


How does Autoritas do it?

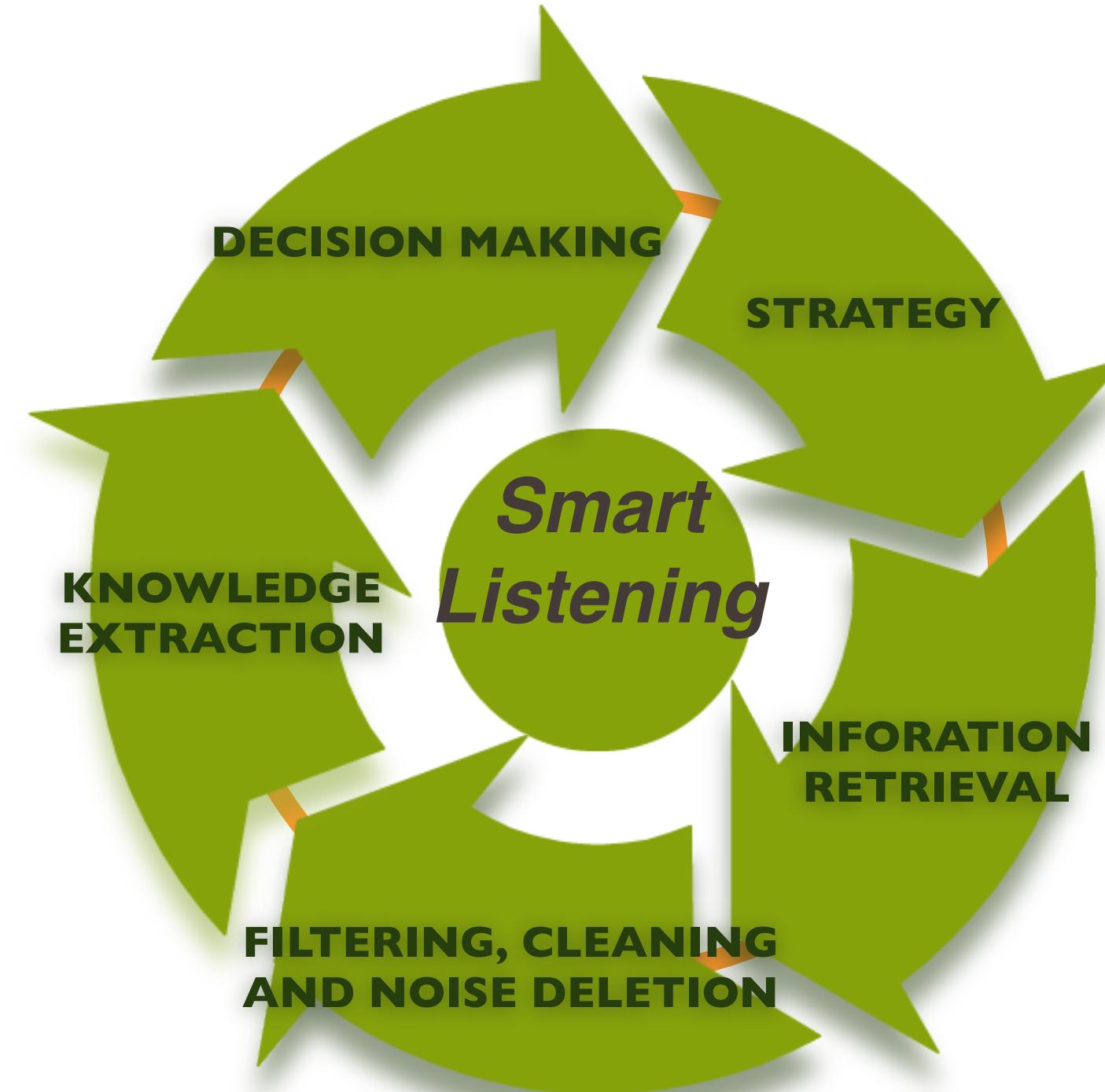
Method

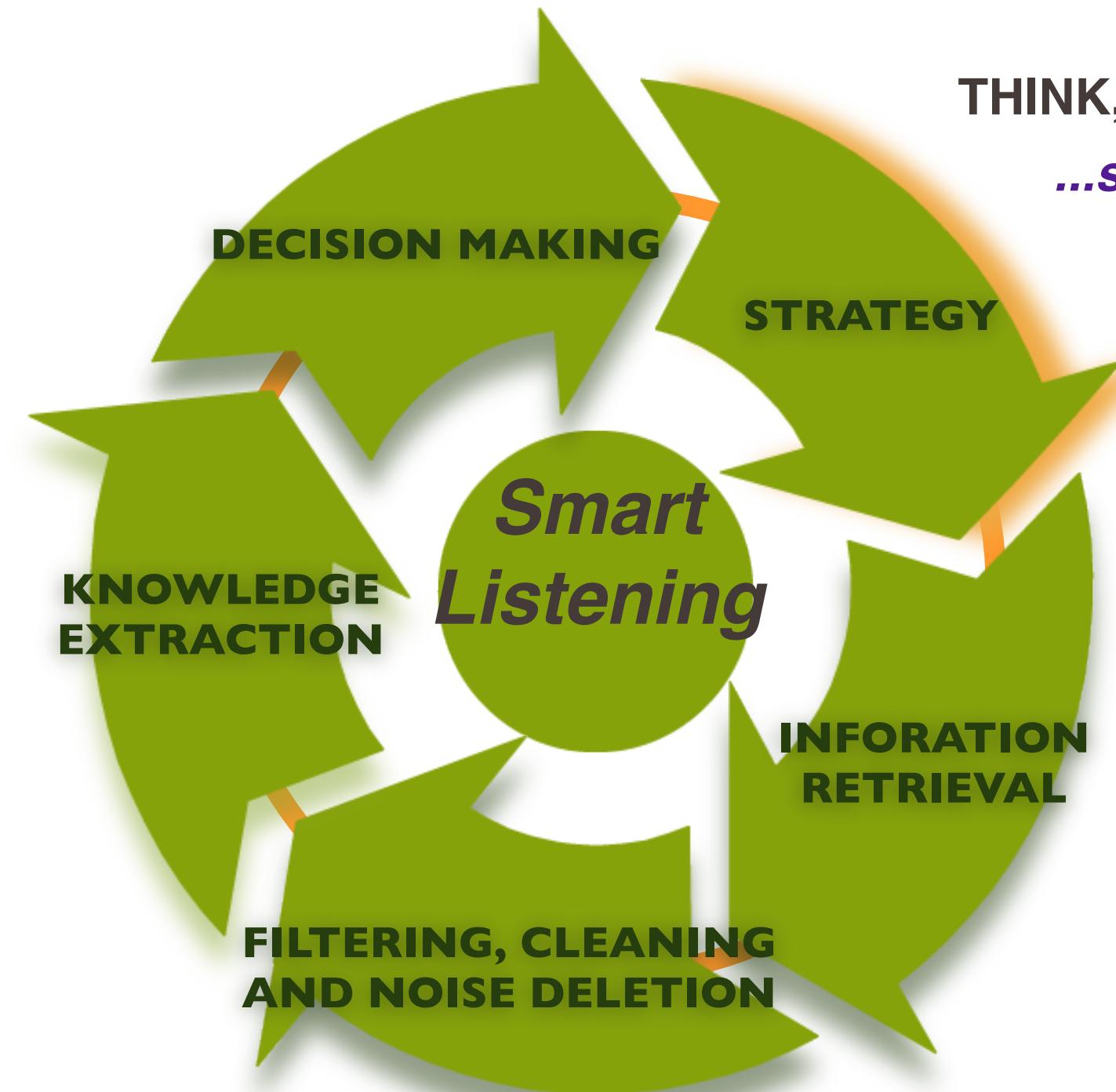


Team



The Smart Cycle





THINK, THINK, THINK...

*...since thought precedes
action*

*What do you expect
from Internet / Social
Networks for you
Business
Objectives?*

*Where is your
ROI?*

Police Investigation Bureau

*They have it
very clear!*



***Catching the
bad***

Tourist Agency

*Demand
prediction*



Insurance Companies

*Measuring Reputation
in risk management*

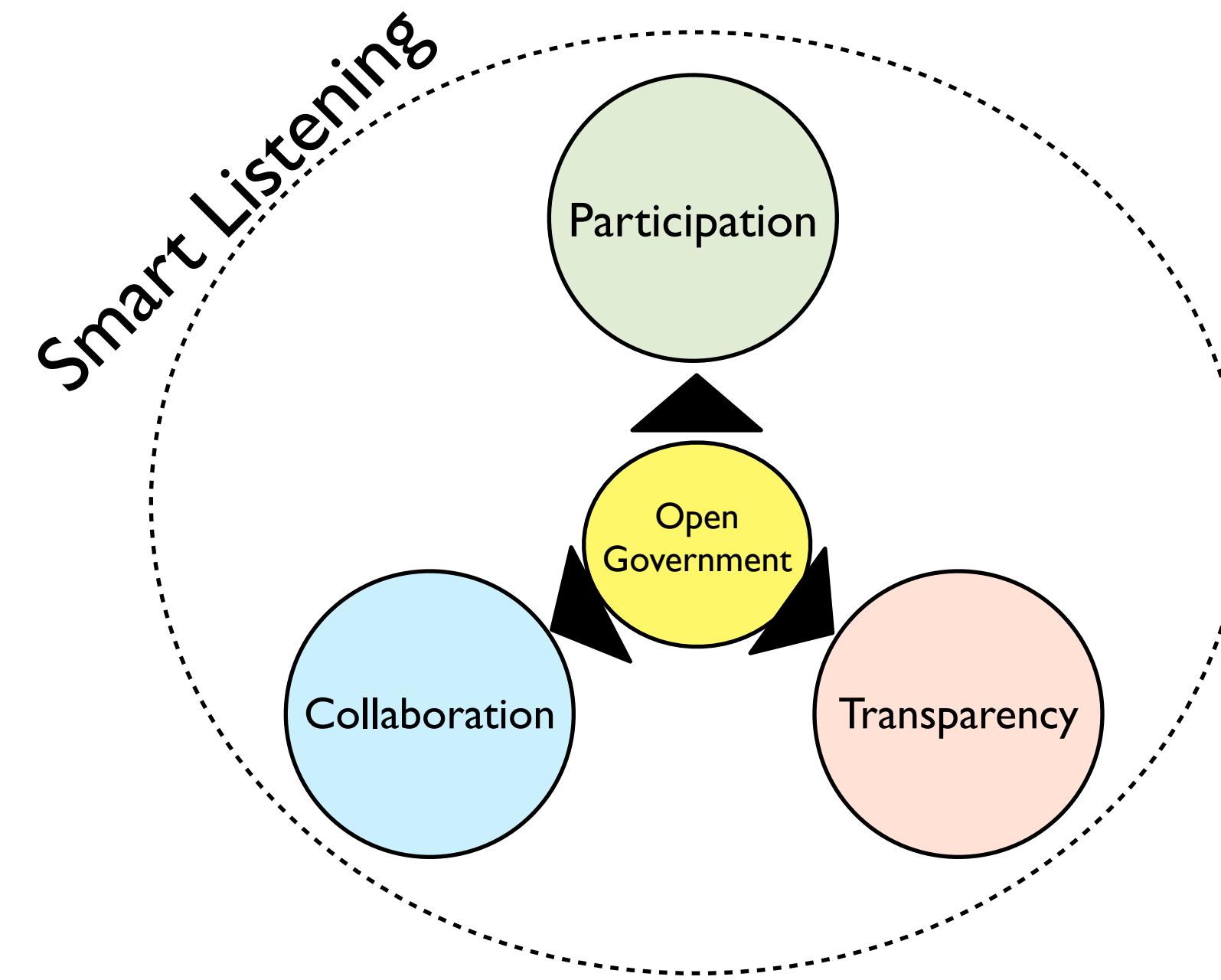


Marketing Agencies

*Internet as
Information source*



Governments



Telcos

*To improve
customer
service*



Mass media



Social content generation

Audience measurement

And yours?

Think, think, think...



“Truth is out there”
X-Files

$$\lim_{x \rightarrow \infty}$$

OBJECTIVE:

To retrieve **all** that we
should retrieve **and/but**
without retrieving **nothing**
that we should not
retrieve

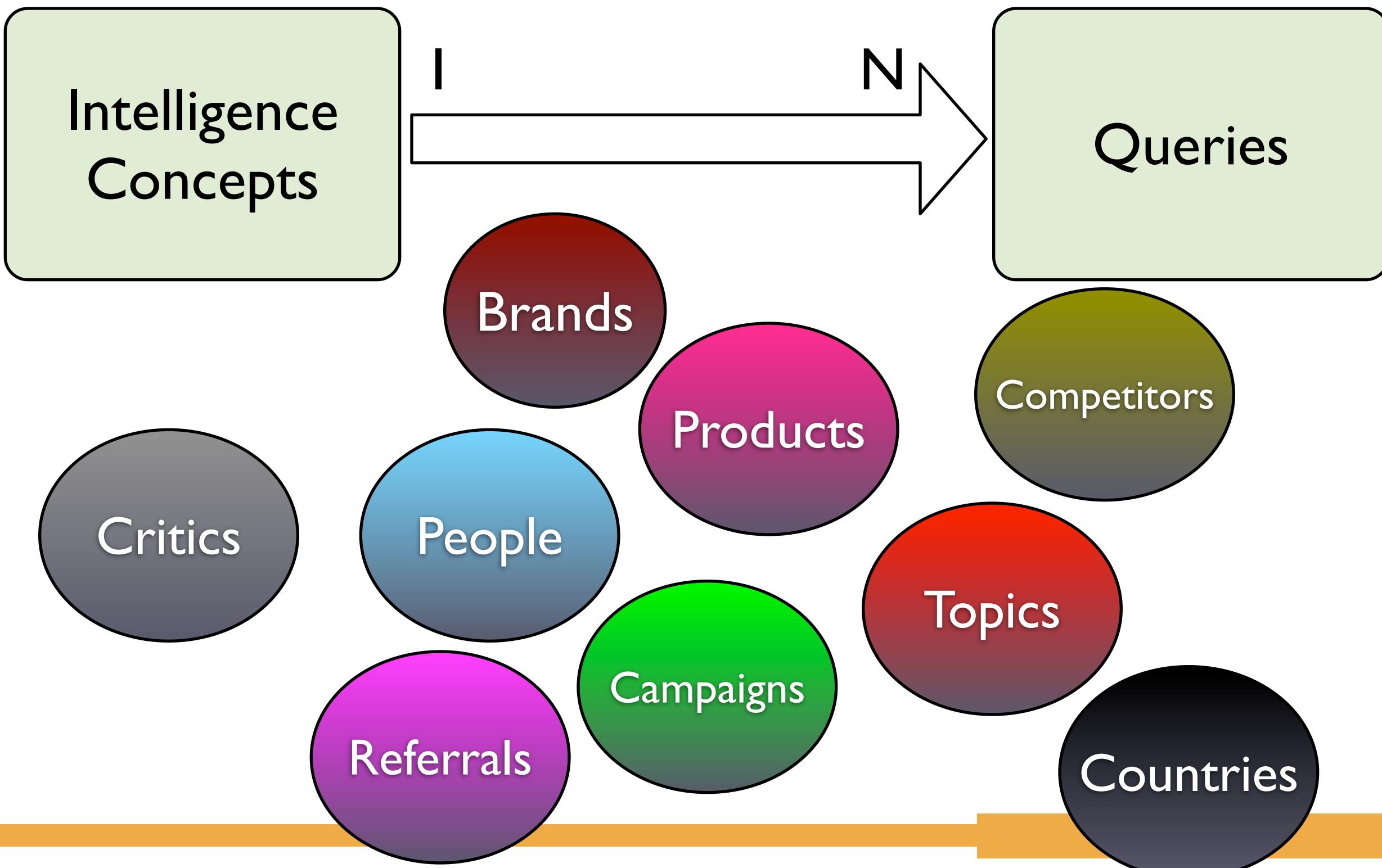
$$\lim_{x \rightarrow 0}$$



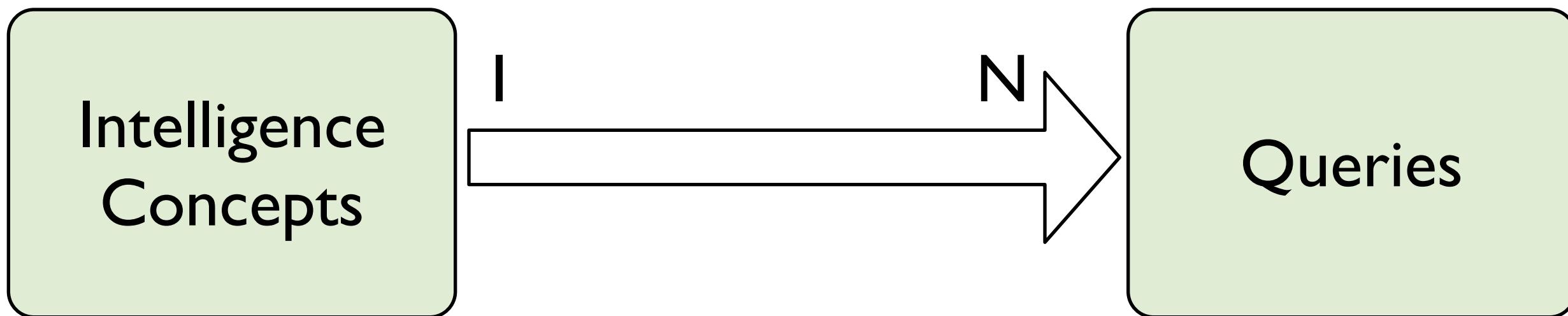
Information Sources (channels)



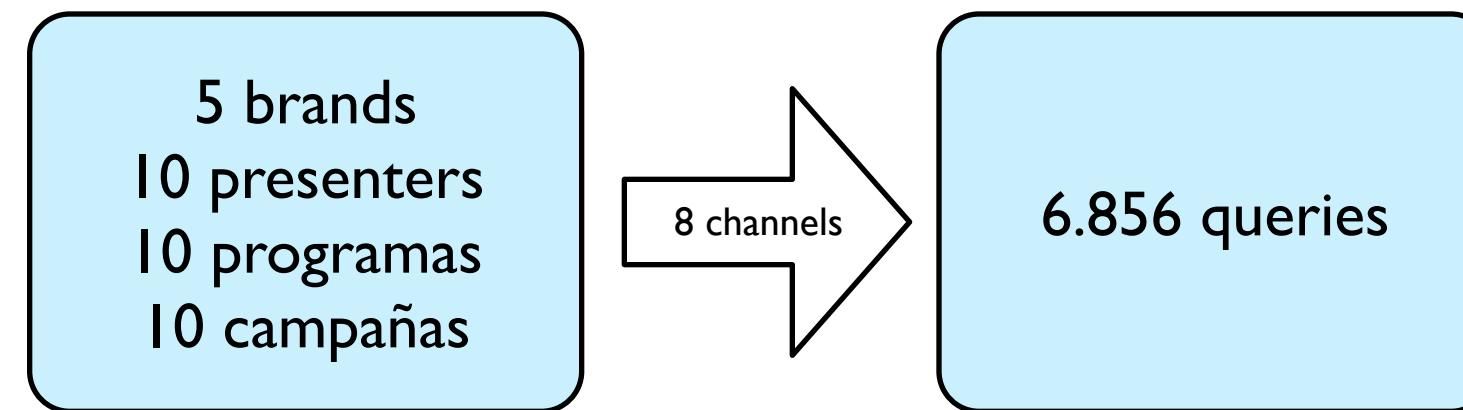
Concept-query mapping



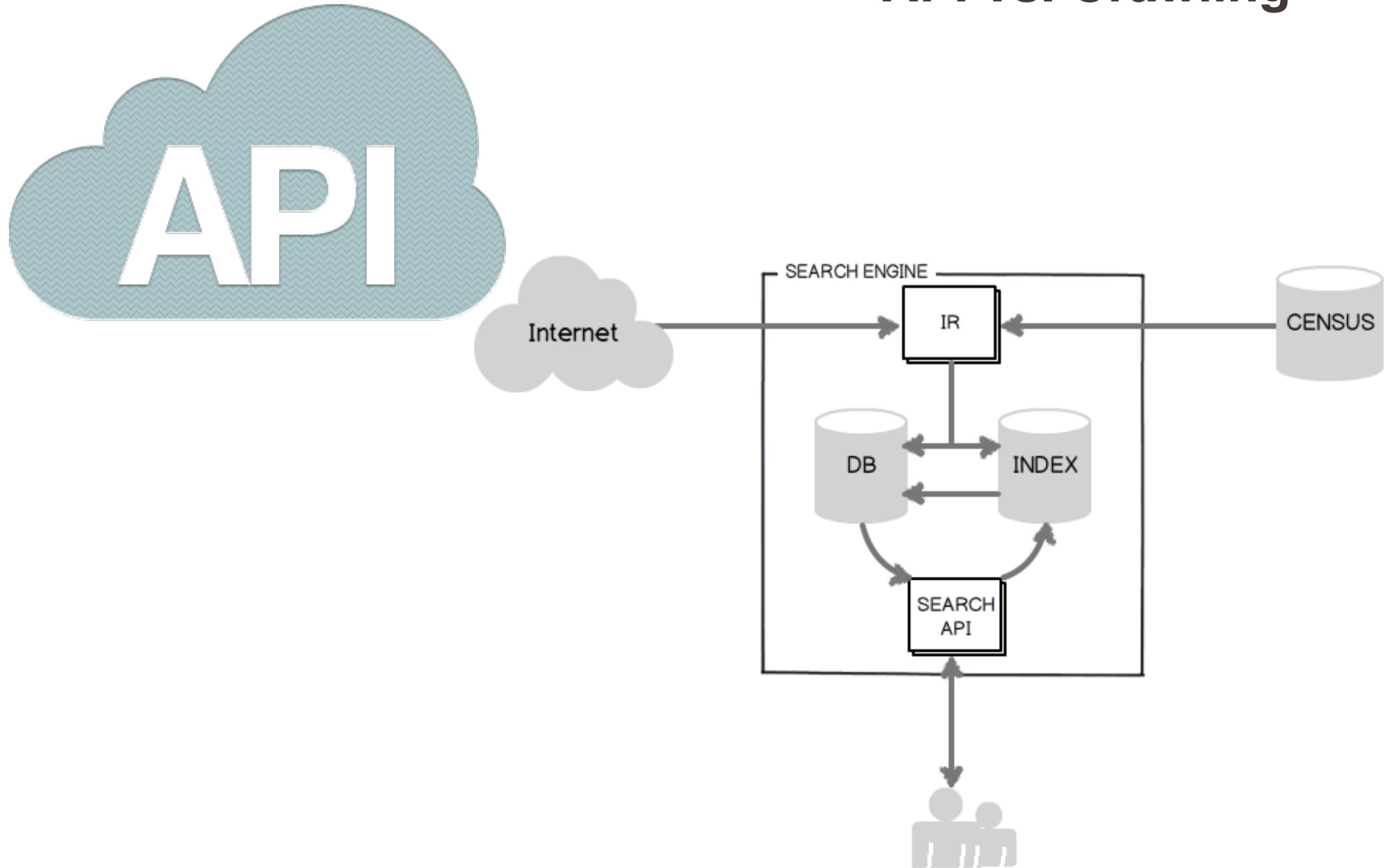
Concept-query mapping



A real example: Spanish public television



API vs. Crawling

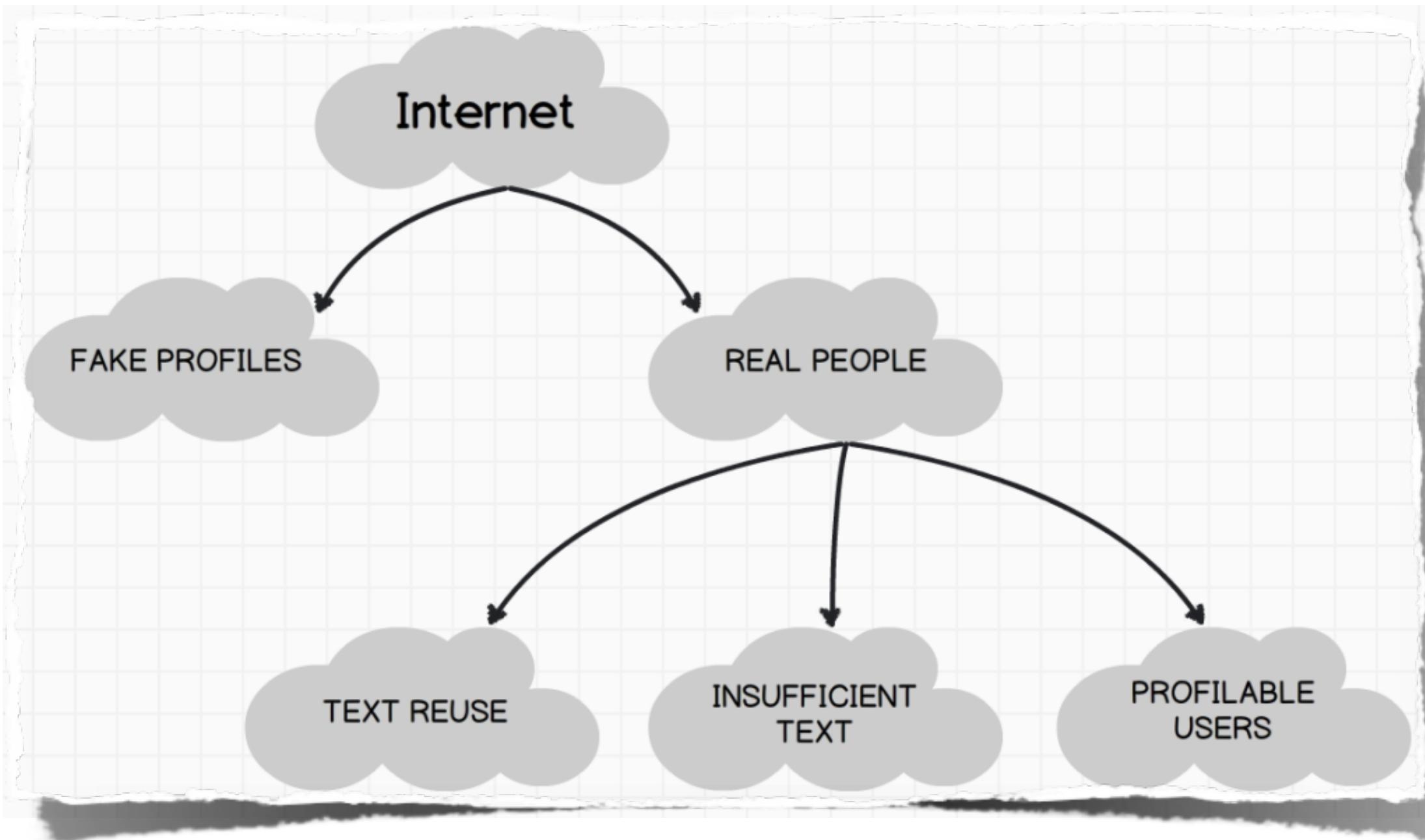




information = data - noise

lim
 $x \rightarrow 0$

Why is there noise?



Users lie, plagiarize or say foolishness...

For example...



findability ≠ *relevance*

TRAGEDIA EN BUENOS AIRES

Casi medio centenar de muertos en un accidente ferroviario en Argentina

- El convoy, que transporta diariamente a más de mil personas, no frenó al entrar en una de las principales estaciones de la capital argentina
- **Video:** Momento en que el tren choca al llegar a la estación Once
- Sigue en directo la transmisión de A24 sobre la tragedia

FRANCISCO PREGIL | Buenos Aires | 22 FEB 2012 - 21:41 CET

Archivado en: Buenos Aires Latinoamérica Argentina Accidentes ferrocarril Sudamérica Américas

Accidentes Sucesos



Rescatistas trasladan a un herido. / JULIO SANDERS (REUTERS)

■ Recomendar
376
■ Twittear
264
■ Enviar

■ Compartir
■ Enviar
■ Imprimir

El tren de cercanías Sarmiento iba con casi todos los viajeros apurados de pie, como siempre en hora punta. Salio a las once y media de la estación de Moreno para recorrer 14 estaciones hasta Buenos Aires. Sobre la estación cuarta, en la de Castelar, cambió de conductor. El nuevo maquinista, de 28 años, iba a emprender su primer trayecto de la mañana. Y el tren siguió frenando y arrancando en cada una de las paradas. Parecía un viaje normal, tal vez un poco más incómodo que otros para sus más de 1.200 viajeros, porque era la primera jornada laborable tras un largo puente de carnaval. A mil metros de su destino redujo la velocidad de 47 a 39 kilómetros por hora. En el andén entró a 26 kilómetros por hora, según el ministro de Transporte de Argentina, Juan Pablo Schiavi. Eran las velocidades normales de entrada en la estación. A 40 metros del final ya había frenado hasta los 20 por hora. Pero ya no volvió a frenar más. De pronto, el tren impactó contra el muro de contención y el segundo vagón se incrustó más de cinco metros en el primero. Eran las 8.32 (las 12.32 hora peninsular española). Murieron al menos 49 personas y 600 resultaron heridas. Uno de los que resultaron con vida fue el propio maquinista de 28 años, quien anoche se encontraba en una unidad de cuidados intensivos. "No sabemos qué ocurrió en los últimos 40 metros", reconoció el ministro.

PUBLICIDAD

AQUÍ Y AHORA

www.altamirasantander.com
902 509 559

Francisco Pregil Pecellin

EL PAÍS

Jazztel ADSL 15,95 €

ADSL por 15,95€ y primer mes gratis. Siempre es mejor estar en buena compañía, elige Jazztel. Con llamadas a móviles gratis. www.jazztel.com

ads by EL PAÍS

E ÚLTIMA HORA

Chile ha reabierto su frontera con Perú después de que la cerrara el lunes por haber localizado 11 minas antipersona de la época de la dictadura de Pinochet. Las recientes lluvias en la zona han desenterrado los artefactos.

EL PAÍS Hace 28 minutos

fwd @el_pais: La aspirante se presenta. Pilar Sánchez Acera presenta su candidatura a la secretaría general del PSM: "Es el momento del cambio". Crónica de José Marcos <http://cort.as/1dJ8>

EL PAÍS Hace 33 minutos

HP cae más de un 2% tras presentar resultados. El fabricante de ordenadores registró una caída del 44% en el beneficio trimestral, hasta los 1.470 millones de dólares. Los ingresos lo hicieron un 7%, hasta 30.040 millones. Las ventas de PC bajaron un 18%, informa Sandro Pozzi.

EL PAÍS Hace 39 minutos

El de hoy, con 49 muertos, es el **tercer accidente ferroviario más grave en la historia de Argentina**, según cuenta el diario La Nación. Está cerca del segundo peor, que dejó 55 muertos y que, como el más trágico (142 fallecidos), sucedió en la década de los 70 <http://cort.as/1dk>

estup

Advertisements without relevance for the contents

Usable contents

Latest news section that distorts the semantics of the page

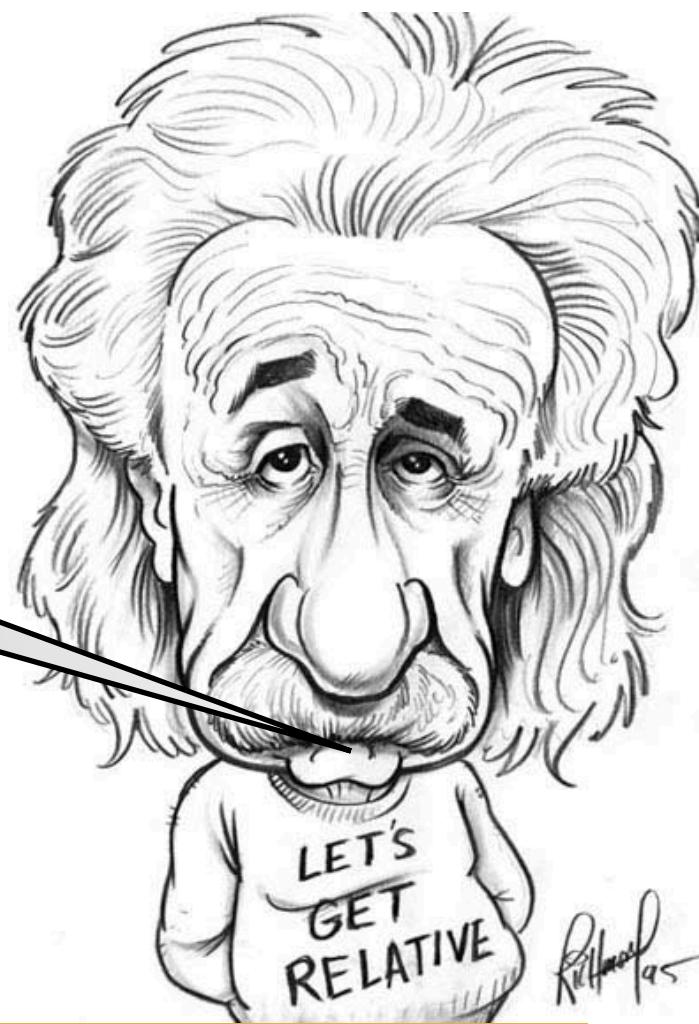
lím
 $x \rightarrow 0$

The importance (and difficulty) of obtaining the right date



If the url includes the date,
it's easy to know it

That's relative. Is
this url from july or
january:
[http://xxx/07/01/2010/
crawler-403-
forbidden.html](http://xxx/07/01/2010/crawler-403-forbidden.html)



lim
 $x \rightarrow 0$



Language filtering

English

estoy sin internet ¬¬ fuuuuck!!!

Finnish

... euskocaja, como euskolabel, euskotren, euskomueble... XDDD

German

Vierrrrrrrrrnes, egunon!!

Portuguese

Flowah Powah!

Language Models vs. n-Gramms vs. Machine Learning

Geography filtering



Koldo M Martin

@iTitanMiller Mi casa, bilbao!

Ufff!!!

nerea miguel andrada

@Nereabskt bilbooo

i lovee baskett!!! BBB & ATHLETIC:)

Jeremy Hagger

@mac_english M.A.R.S.

The Love Jones

<http://twitter.com/lovejones>

lim
 $x \rightarrow 0$

Fernanda

@FernandNavarro Narnia

Exactly who we are is just enough.

Source geography vs. contents geography vs. profile geography

“Limpia, fija y da explendor”

Lema de la RAE

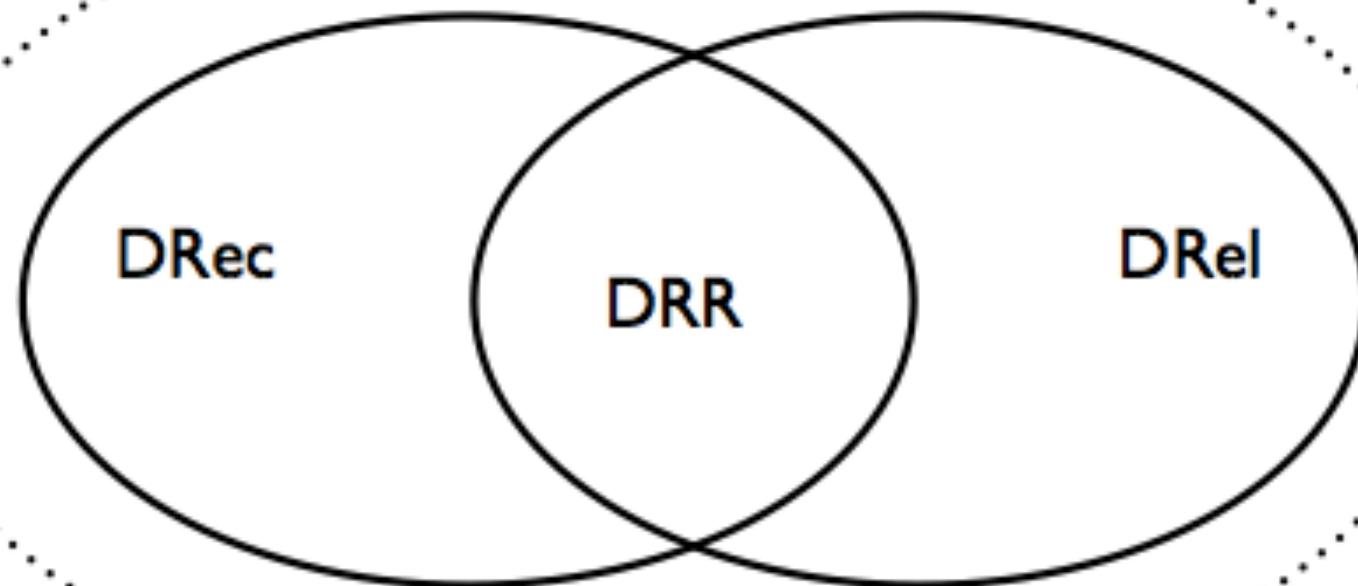
80% of
the
workload



I believe I
solved the
equation:
 $\lim_{x \rightarrow \infty} \lim_{x \rightarrow 0}$



Information Retrieval Evaluation...



DRel: Documentos Relevantes

DRec: Documentos Recuperados

DRR: Documentos Relevantes Recuperados

$$\text{Precisión} = \text{DRR} / \text{DRec}$$

$$\text{Alcance} = \text{DRR} / \text{DRel}$$

- *Do you want that all retrieved content is good? -> Precision*

$\lim_{x \rightarrow 0}$

- *Don't you want to leave anything without retrieving? -> Recall*

$\lim_{x \rightarrow \infty}$

...in the Science

Information Retrieval Evaluation...

7.000 retrieved
54 wrong
99.23% precision



3.000 retrieved
50 not retrieved
98.36% recall

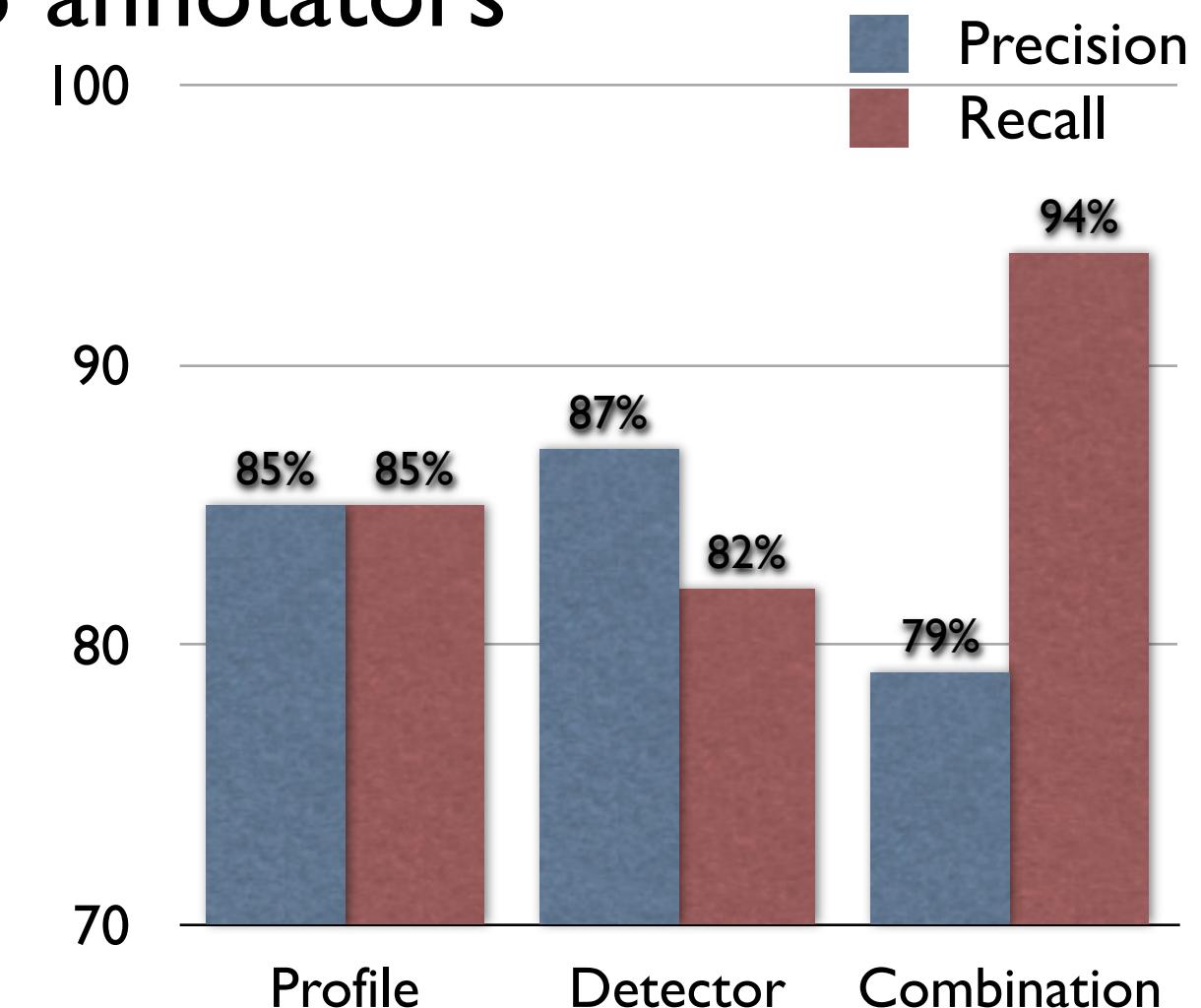
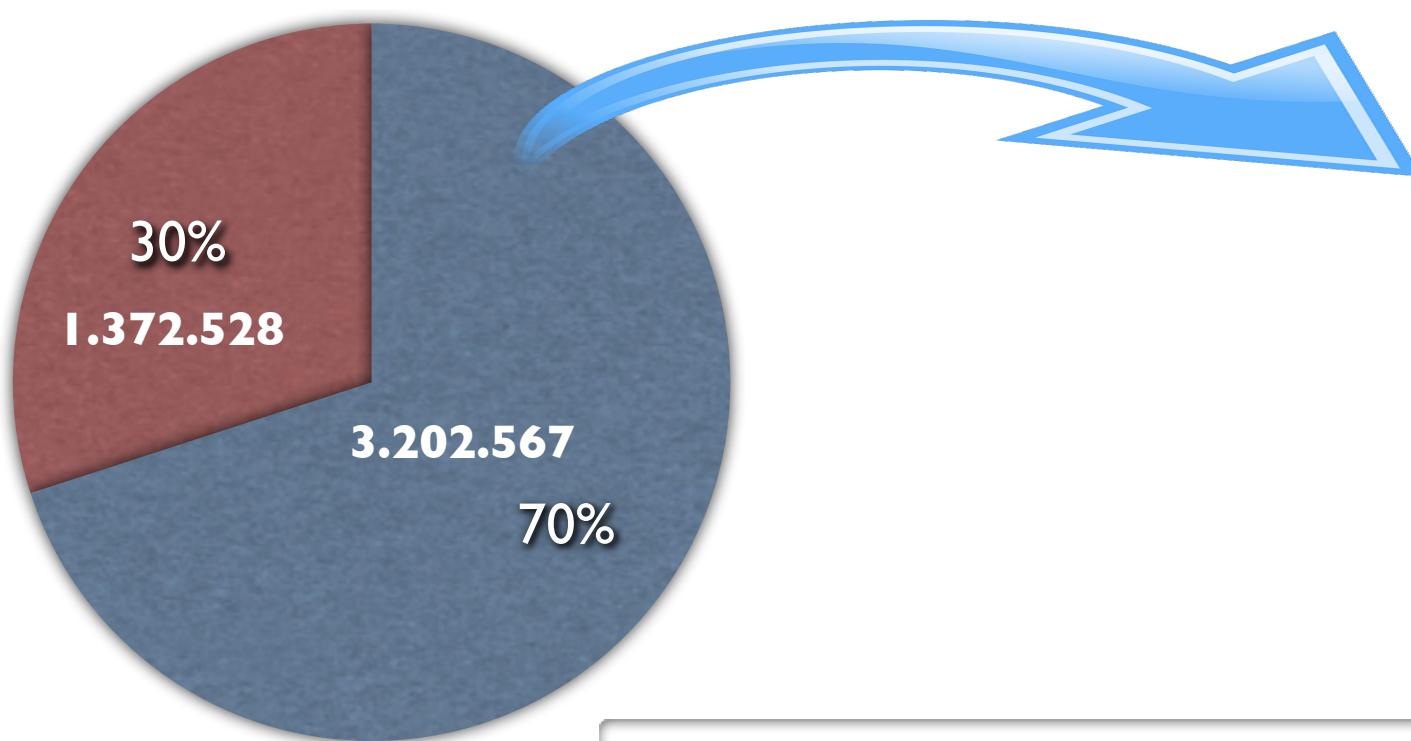
...in the Industry

Touristic project: 4.575.096 tweets

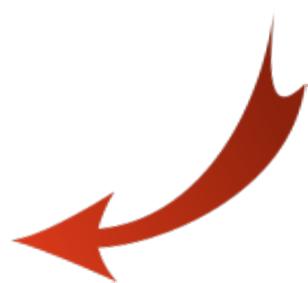
QUERIES: Mallorca, Menorca, Ibiza, Formentera, Baleares, ...

METHOD: 1.200 random tweets / 3 annotators

- Labeled
- Not labeled



672.539 wrongly retrieved
204.419 good but not retrieved
1.372.528 that we don't know...



Bring order to the information...

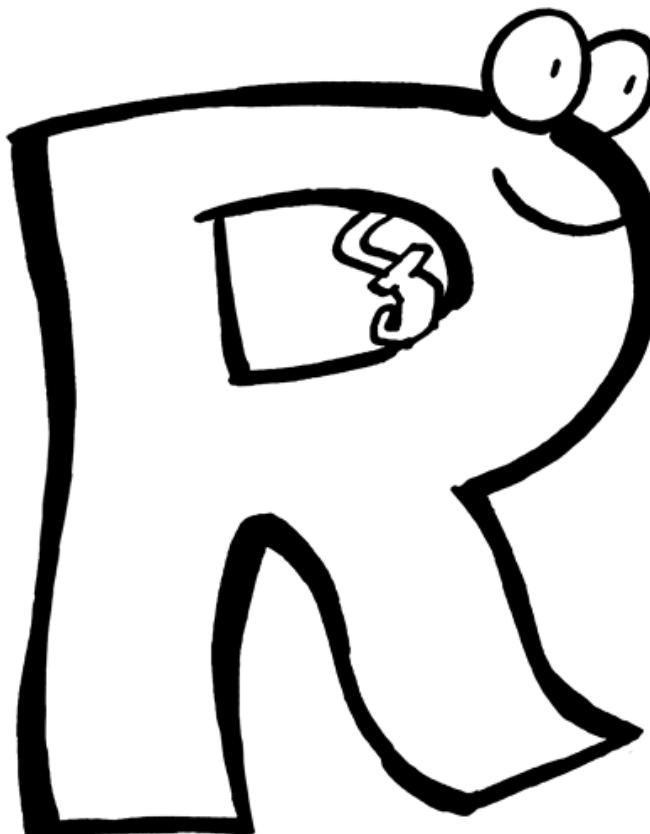
...to analyse it and get value



...to ask questions...



...and to know which new questions to ask



**WHAT are people talking about
the consonant, the prefix,
language or the company?
telecommunications?**



Semantic ambiguity in millions of documents!!!

WHEN? -> Crisis management



When are things happening?

WHERE/FROM WHERE is people talking about it?

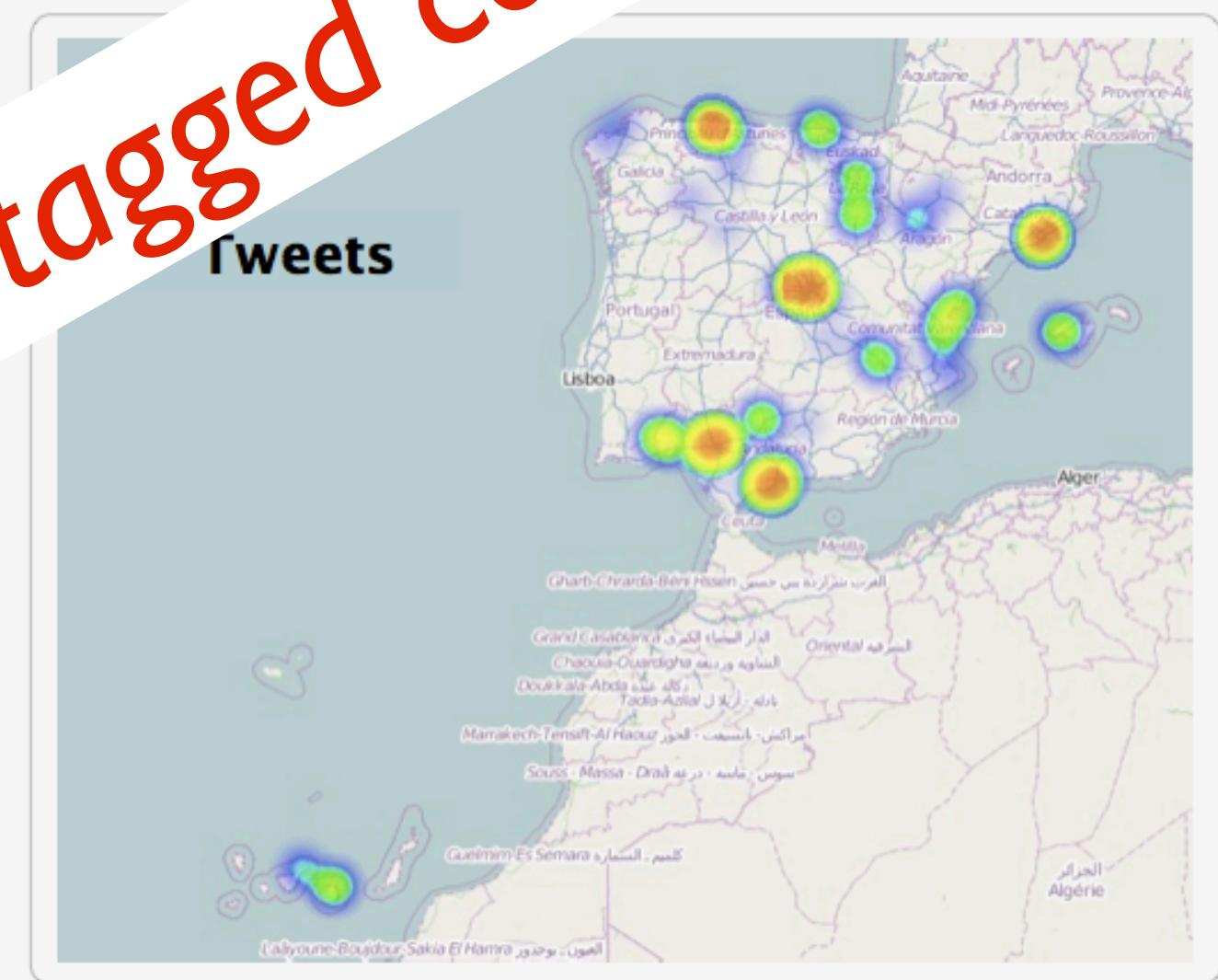
Distribución geográfica

Votos



Mapa de calor en función del reparto de votos por mesas electorales

Tweets



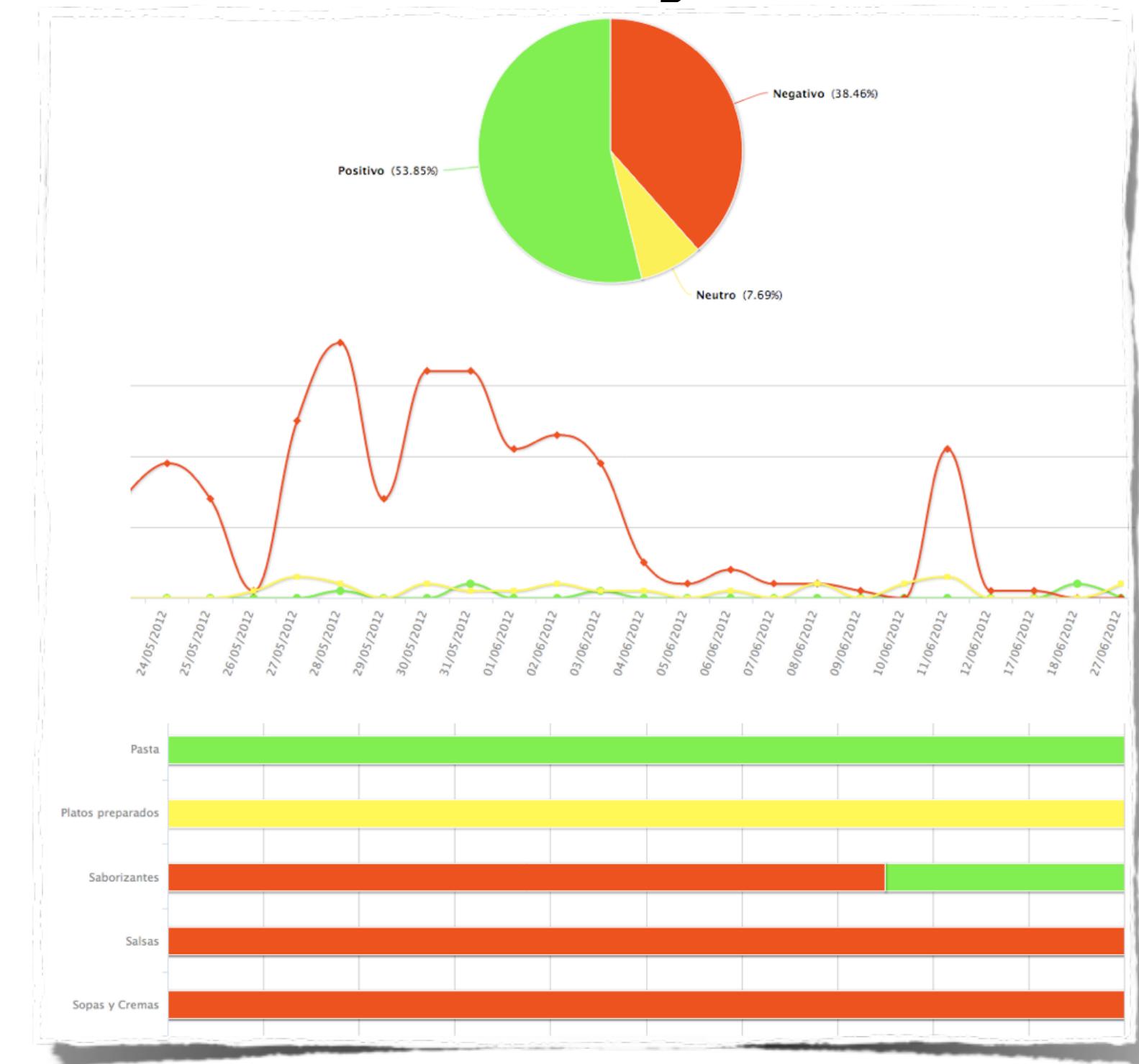
Mapa de calor de los tweets emitidos, que han podido geolocalizarse

HOW? -> Not only sentiment analysis

Polarity is onlye one dimension:

- Emotions
- Motivations
- Values
- SWOT

All of them answer the question “how?”



An example: “The risk premium in Spain is 235”

Positive, negative, neutral o none?

An example: “The risk premium in Spain”

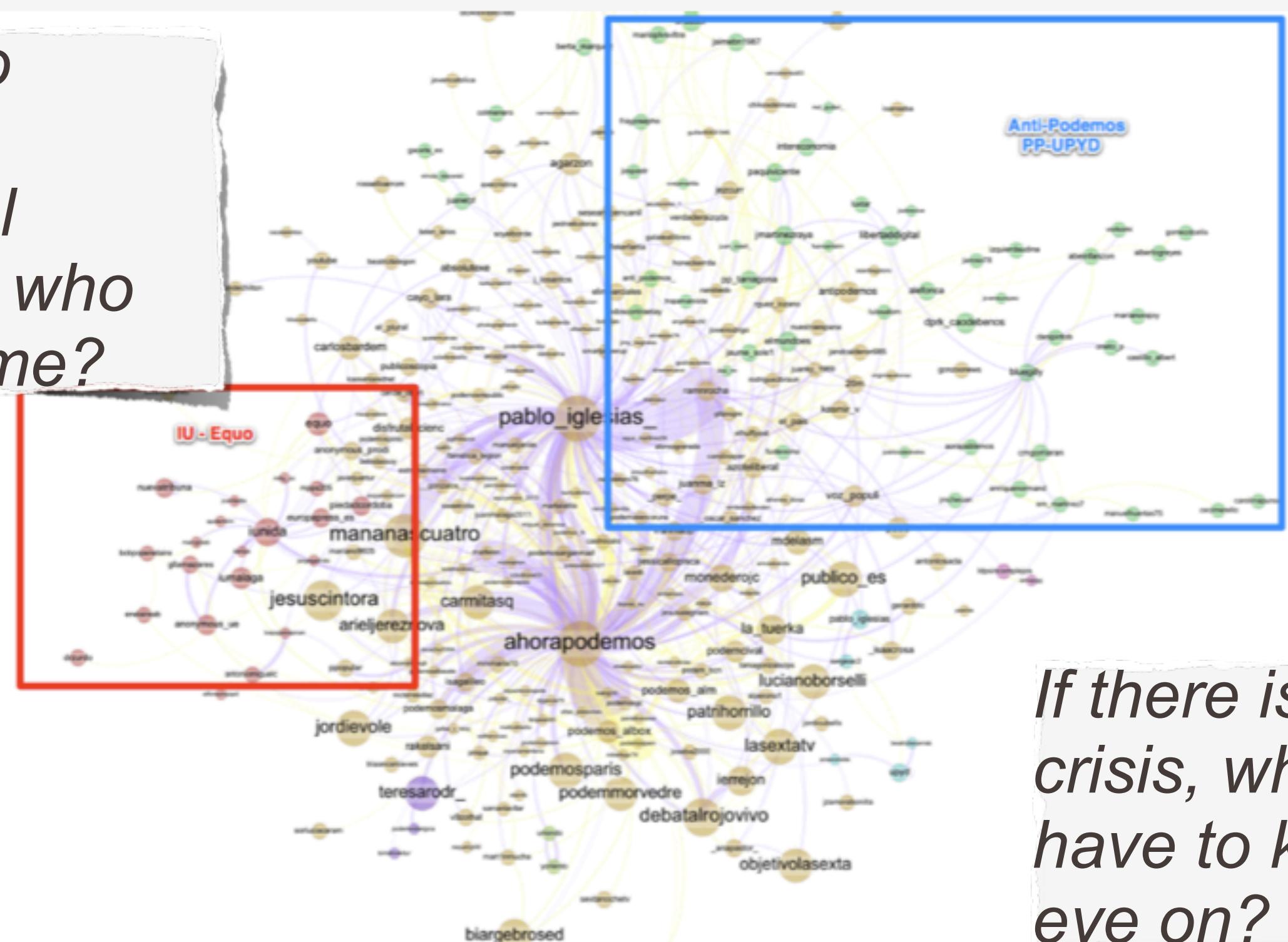
Positive, negative, neutral?

My question: For

Subjectivity in transmitter,
and in receiver!!
country?
opposition?
. of the Spanish Bank?
foreign investor?
national capitalist?
or whom has a mortgage?

WHO? -> Social Network Analysis

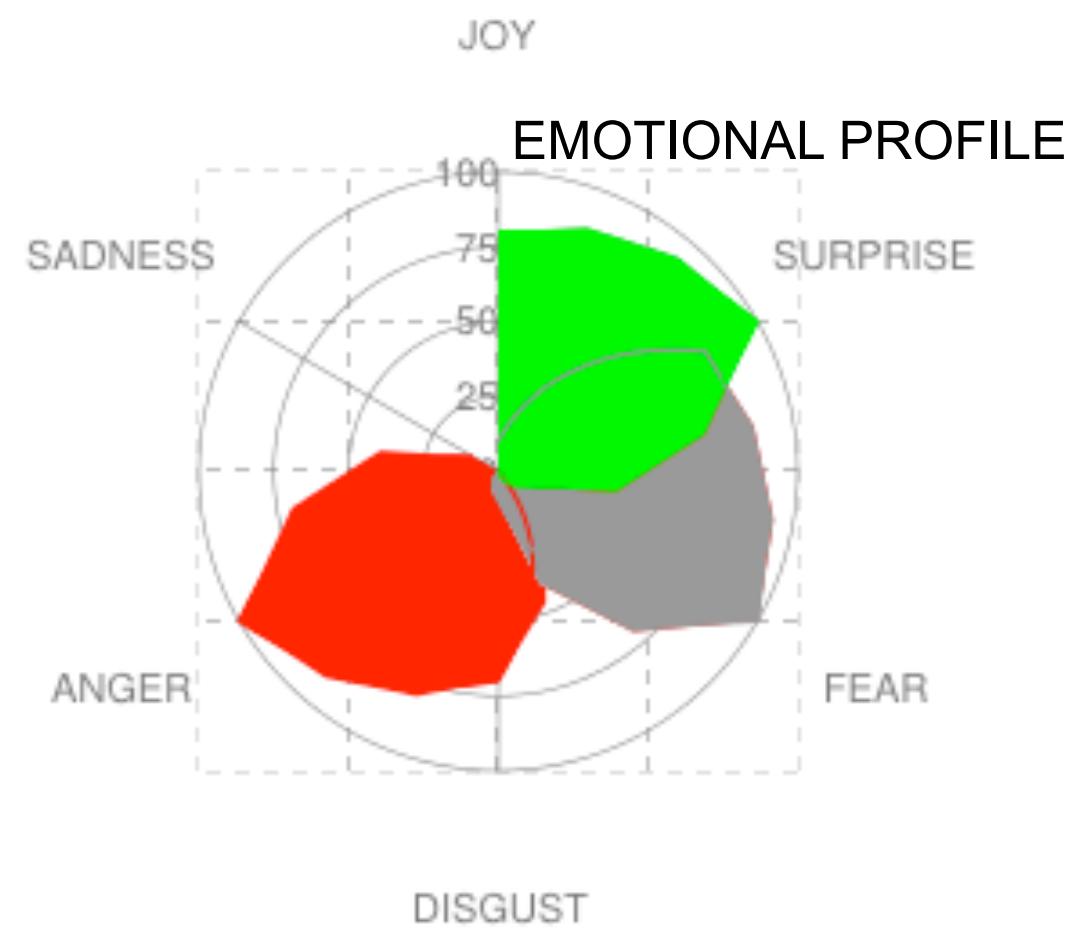
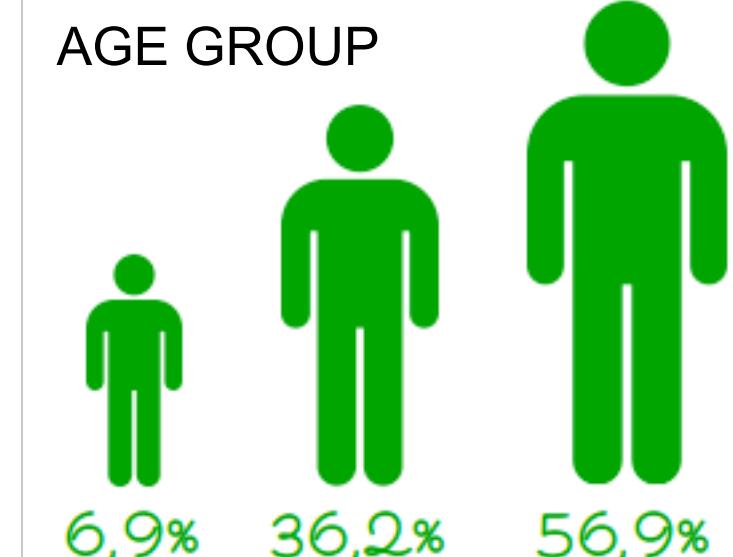
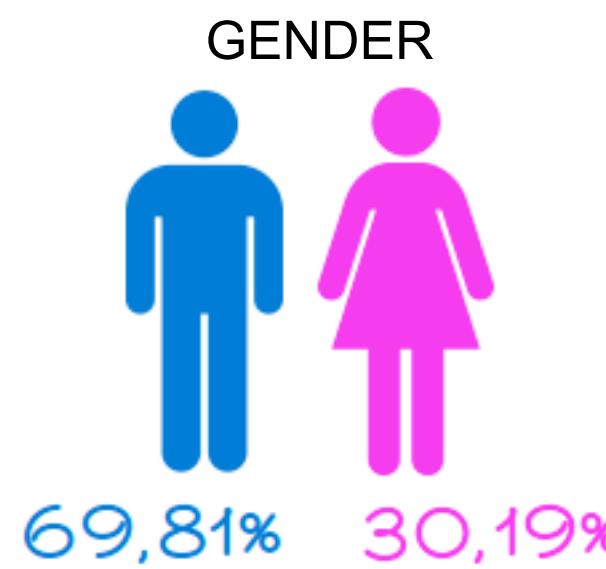
If I want to transmit a successful message, who can help me?



Usuarios de twitter más influyentes en la conversación sobre PODEMOS

If there is a crisis, who do I have to keep an eye on?

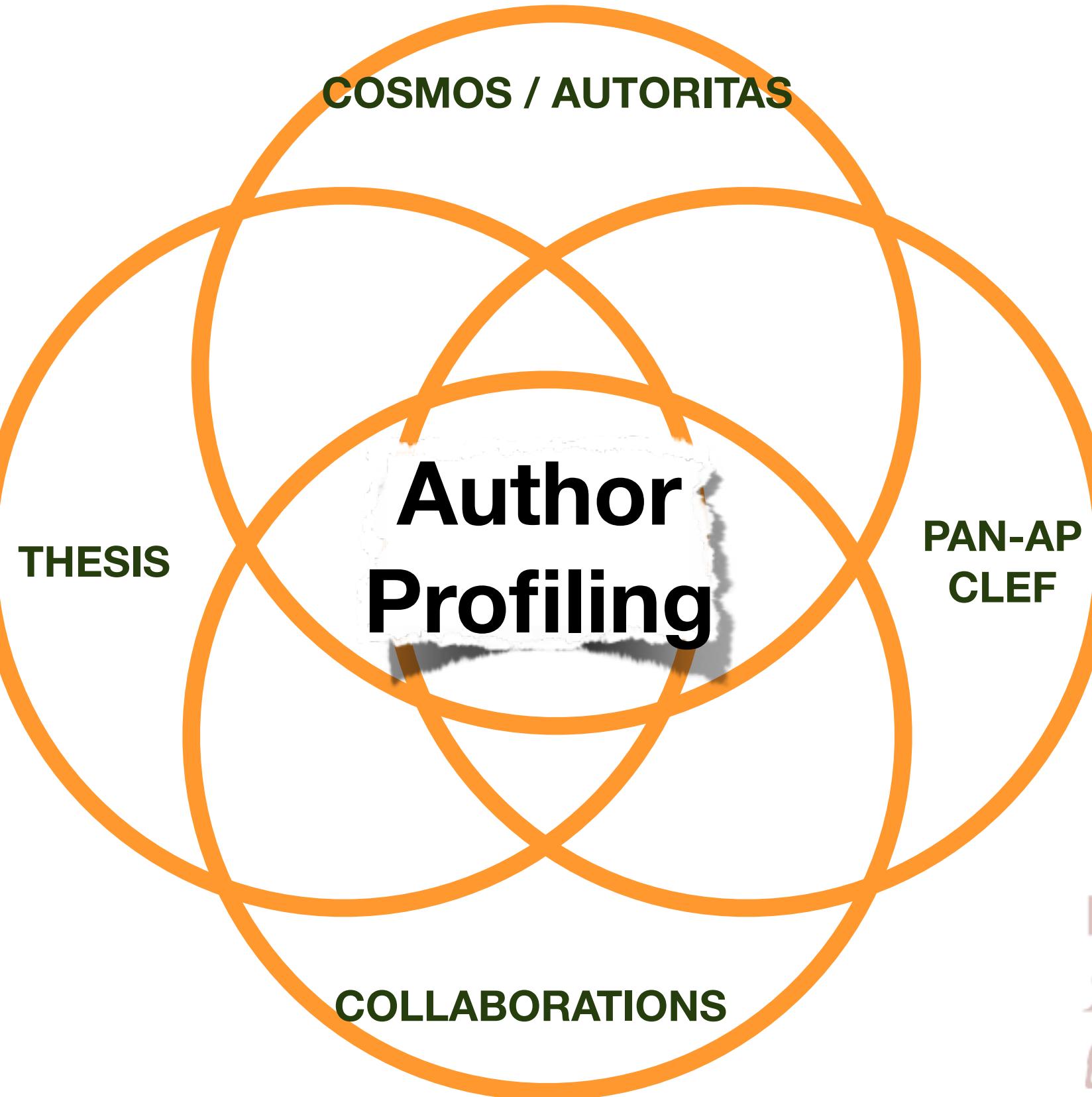
WHY -> Author Profiling



PERSONALITY TRAITS



... political ideology, religious beliefs, and much more!



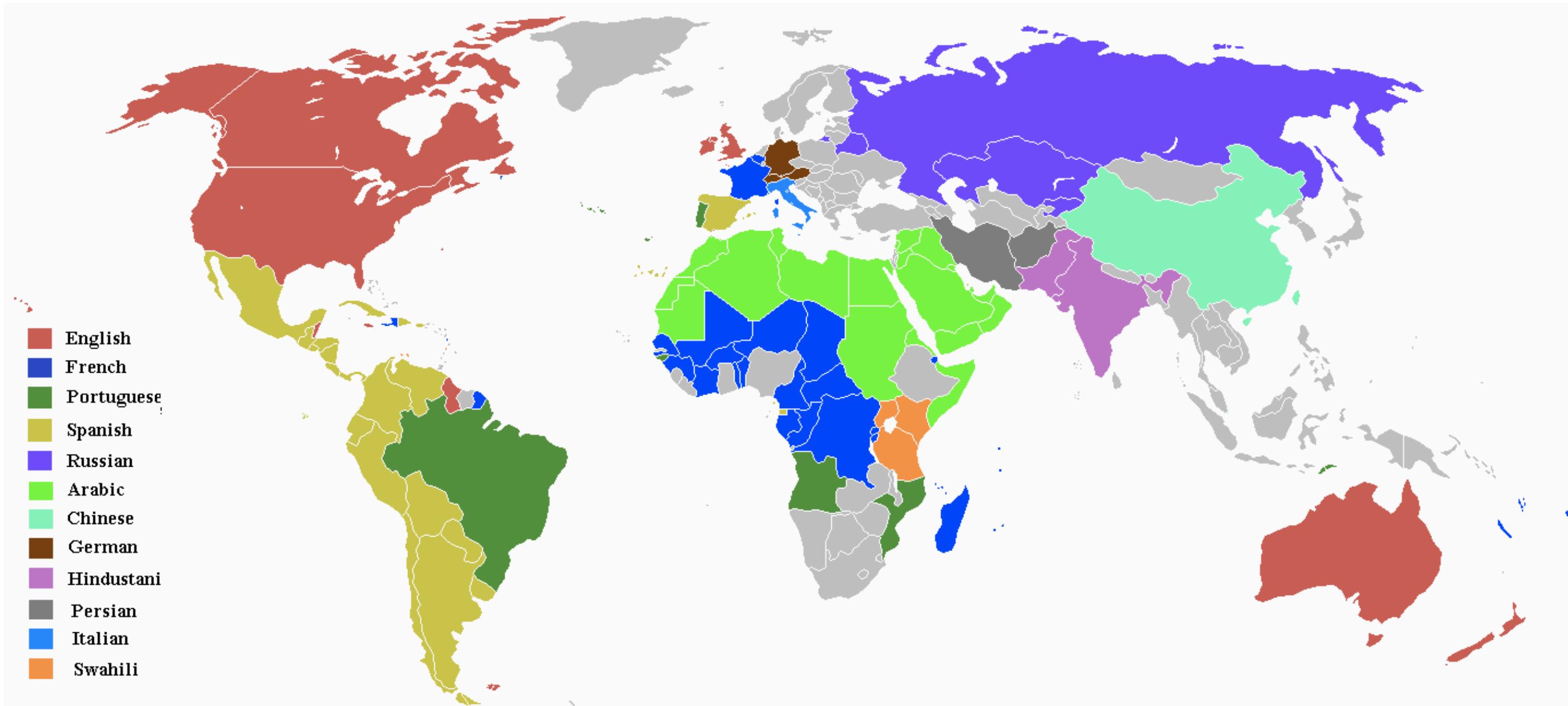
In collaboration with:
Paolo Rosso



UNIVERSITAT
POLITÈCNICA
DE VALÈNCIA

Language Langue Linguaggio
Языка **NLEL** SPRACHE
Lingua LLENGUATGE **NLEL**
Natural Language Engineering Lab

On the Internet there are no boundaries...



...except the language

Show me how you talk...



HispaBlogs

TRAINING	TEST
450	200
x 5 varieties	



...and I tell you where you are

Automatic Identification of Language Varieties: The Case of Portuguese.
Zampieri, M., Gebrekidan-Gebre, B.
In Proceedings of the Conference on Natural Language Processing 2012

- * Corpus (1 000 documents from newsletters): 2 regional variations
- * Features: word and character n-grams
- * ML Algorithm: Language probability distributions with log-likelihood function for probability estimation
- * Evaluation method: 50/50 split
- * Accuracy:
 - * Word uni-grams: 99.6%
 - * Word bi-grams: 91.2%
 - * Character 4-grams: 99.8%

Automatic Identification of Arabic Language Varieties and Dialects in Social Media.
Sadat, F., Kazemi, F., Farzindar, A.

In Proceeding of the 1st. International Workshop on Social Media Retrieval and Analysis SoMeRa 2014

- * Corpus (blogs and forum documents): 6 regional variations
- * Features: character n-grams
- * ML Algorithm: Markov language model vs. Naïve Bayes
- * Evaluation method: 50/50 split
- * Accuracy: 98% (78% F-measure)

Our approach: A probabilistic framework

- * Step 1: Vector Space Model based on TF/IDF

$$w = \ln(tf + 1) \cdot \ln\left(\frac{N}{1+idf}\right)$$

w: weight of the term
 tf: term frequency
 idf: inverse document frequency
 N: number of documents

- * Step 2: Build the supervised matrix

$$D = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} & c_i \\ w_{21} & w_{22} & \dots & w_{2n} & c_i \\ \dots & \dots & \dots & \dots & \dots \\ w_{m1} & w_{m2} & \dots & w_{mn} & c_i \end{bmatrix}$$

W_{mn}: weight of the term n in the document m
 C_i: class of the document i

- * Step 3: For each class in C

$$P(t_j/c_i) = \frac{\sum w_{ij}/c_i}{\sum w}$$

P(t_j/c_i): probability of belonging of the term j to the corpus i
 w_{ij}: weight if the term j in the document i
 c_i: class i
 w: sum of weights

- * Step 4: Calculate class probabilities as the average probability of the existing terms

$$P(c_i) = \frac{\sum P(t_j/c_i)}{n}$$

P(t_j/c_i): probability of belonging of the term j to the corpus i
 n: number of cluster terms in the document

Our approach: document representation

- * A document is represented as:

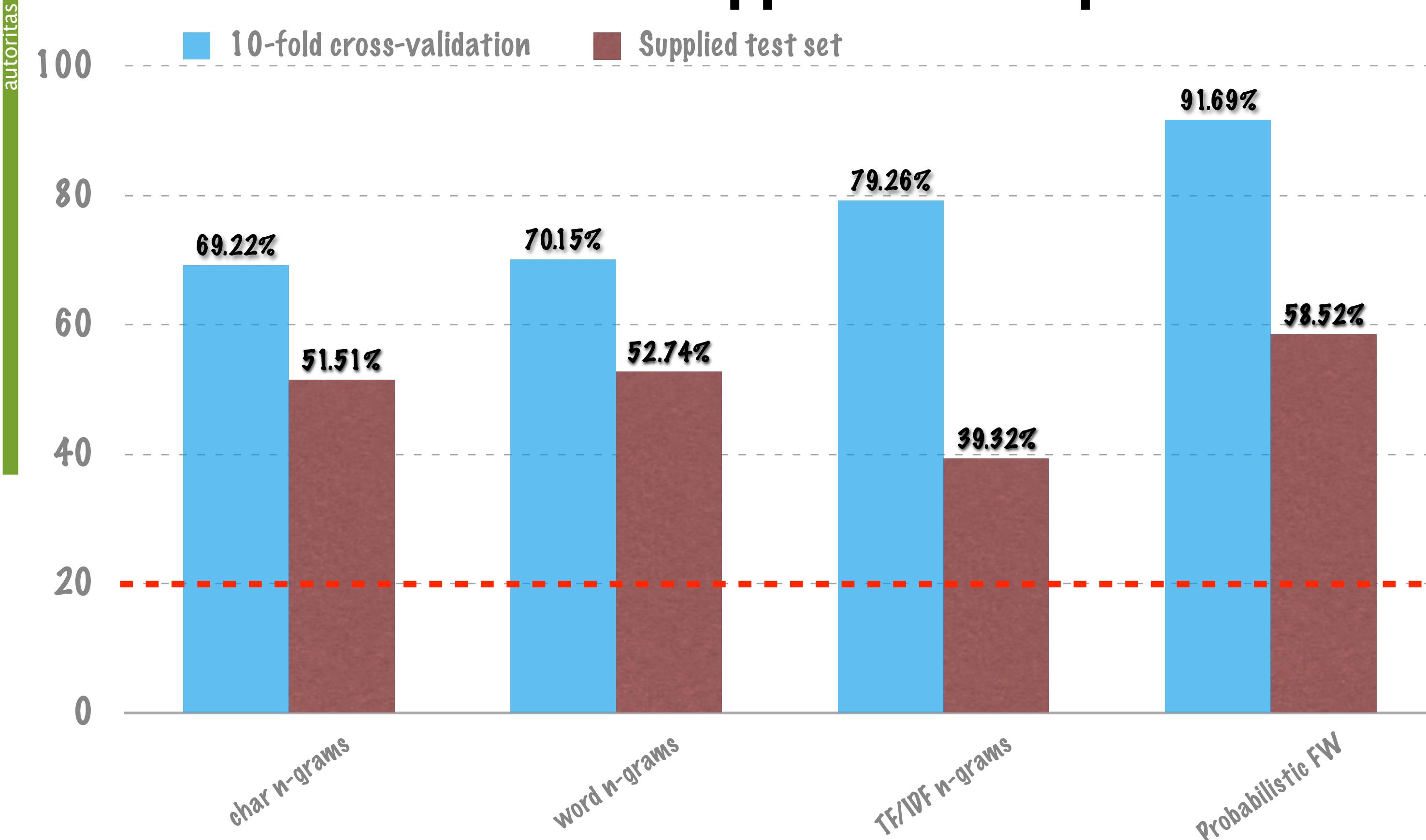
$$d = \{ P(c1), P(c2), \dots, P(cn); ci \}$$

where:

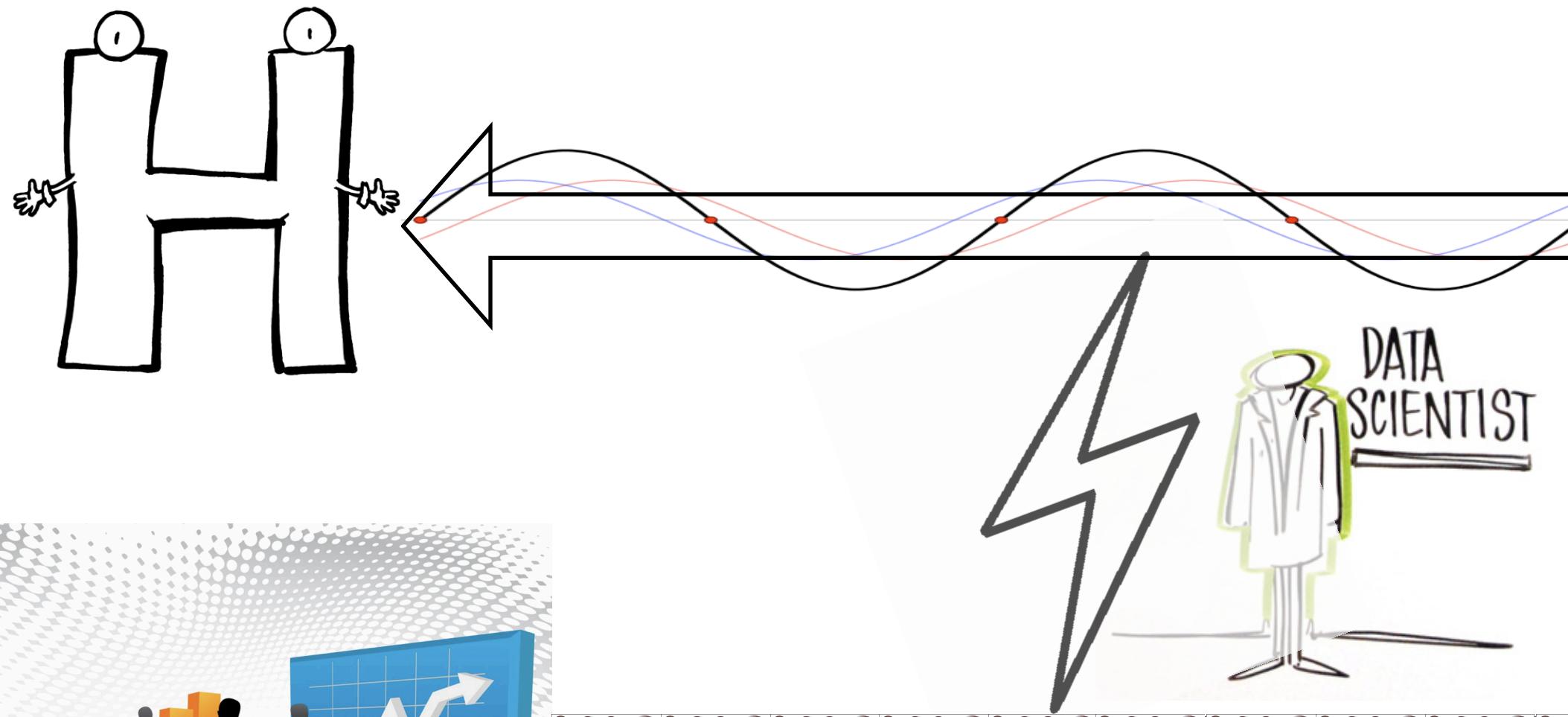
$$P(ci) = \{ \text{avg}, \text{std}, \text{min}, \text{max}, \text{prob}, \text{prop} \}$$

- * There are a high dimensionality reduction: from thousands to 6*number of classes (6*5=30 in our case)

Our approach: compared results



Science vs. Industry



Big Data is both the solution and the problem... ...but mainly the opportunity



- ▶ Retrieval & Storing
- ▶ Evolution
- ▶ Words & Topics
- ▶ Tagging
- ▶ Hashtags
- ▶ People
- ▶ Locations
- ▶ Brands
- ▶ Sentiment & Emotion
- ▶ Users & Relationships
- ▶ Influence
- ▶ Gender & Age
- ▶ Language Variety
- ▶ ...

80~120
tw/sec

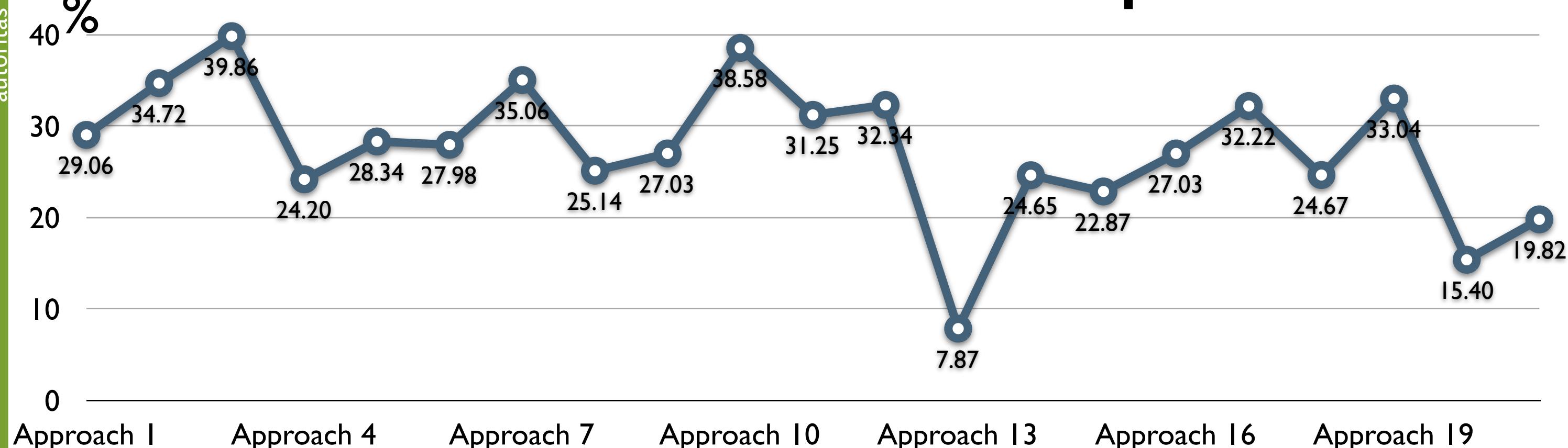
=4.800~7.200
tw/min

=288.000~432.000
tw/hour

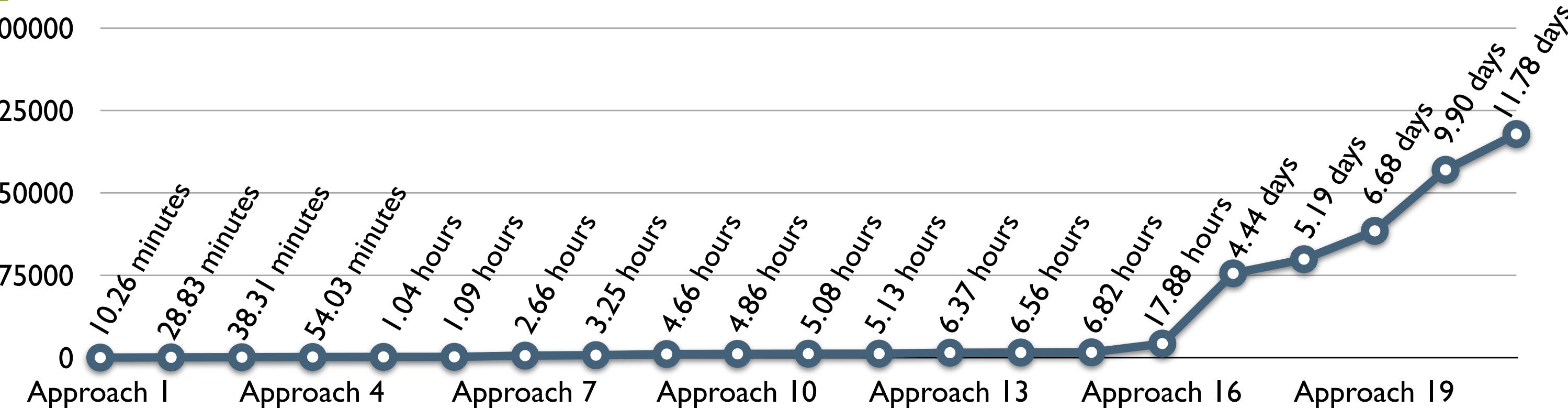
=6.912.000~10.368.000
tw/day

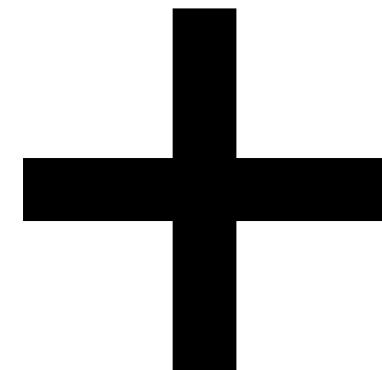
Precision vs. computational cost

consulting, s.a.



AGE AND GENDER IDENTIFICATION ~240K AUTHORS





Technical Skills

“Non-technical” Skills



tocamos**internet**