# Introducción al Procesamiento del Lenguaje Natural | ProgPLN

Víctor Peinado v.peinado@filol.ucm.es 3 de octubre de 2014

¿Qué es el PLN?

El Procesamiento del Lenguaje Natural (PLN) es el estudio científico del lenguaje desde un punto de vista computacional.

Es un área claramente multidisciplinar: lingüística, ingeniería, inteligencia artificial, informática, psicología, etc.

El PLN se interesa en proporcional modelos computacionales de distintos fenómenos lingüísticos. Estos modelos pueden tener dos aproximaciones diferentes:

- 1. sistemas basados en conocimiento: en problemas que podemos modelar, proporcionamos conocimiento lingüístico formalizados
- 2. sistemas basados en estadística: en problemas que no podemos modelar, proporcionamos ingentes cantidades de datos (colecciones de documentos) y dejamos que la máquina cree el modelo a partir del cálculo de probabilidades y la detección de patrones de uso.

## Tareas típicas del PLN

Una buena manera de conocer los temas que tratan un área de investigación es revisar el programa de los congresos más importantes del área:

- ACL 2014: call for papers<sup>1</sup> y programa<sup>2</sup>
- COLING 2014: call for papers<sup>3</sup> y programa<sup>4</sup>
- SEPLN 2014: *call for papers*<sup>5</sup> y programa<sup>6</sup>

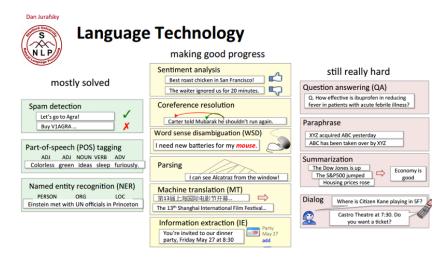
De este modo, podemos identificar algunas de las tareas:

- Desambiguación semántica (word sense disambiguation)
- Análisis morfo-sintáctico (PoS tagging/parsing)
- Traducción automática (machine translation): Google Translate
- Extracción de información (information extraction): TripIt
- Reconocimiento del habla (automatic speech reconition) y síntesis de voz (speech synthesis): Google Voice Search

- <sup>1</sup> http://www.acl2012.org/program/sub03.asp
- <sup>2</sup> http://acl2014.org/Program.htm
- 3 http://www.coling-2014.org/call-forpapers.php
- 4 http://www.coling-
- 2014.org/schedule.php
- <sup>5</sup> http://www.taln.upf.edu/pages/sepln2014/es/callfor-papers.html
- 6 http://taln.upf.edu/pages/sepln2014/es/program.html

- Recuperación de información (information retrieval): Google Search, Bing y Wolfram | Alpha
- Resumen automático (automatic summarization)
- Búsqueda de respuestas (question answering): Ask.com, Watson
- Análisis de opiniones (sentiment analysis) NaturalOpinions
- Comprensión del lenguaje natural (natural language understanding): Siri y Ok Google

Problemas resueltos y cuestiones abiertas



# ¿Por qué es tan difícil el PLN?

El lenguaje natural es eminentemente ambiguo: es la principal diferencia entre lenguas naturales y lenguajes artificiales.

Esta ambigüedad existe a varios niveles:

- ambigüedad fonética y fonológica: vaca/baca, casa/caza, has sido tú/has ido tú\*
- ambigüedad morfológica: casa, beso, río
- ambigüedad sintáctica: Ayer me encontré a tu padre corriendo
- ambigüedad semántica: banco, pie,
- ambigüedad de discurso: correferencia, resolución de anáforas

#### non-standard English Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself

## segmentation issues

the New York-New Haven Railroad

the New York-New Haven Railroad

# idioms

#### dark horse get cold feet lose face throw in the towel

# should never give up either♥ neologisms

unfriend Retweet bromance

### world knowledge

Mary and Sue are sisters. Mary and Sue are mothers.

### tricky entity names

Where is A Bug's Life playing. Let It Be was recorded ... ... a mutation on the for gene

Según la ACL (Association for Computational Linguistics): Computational Linguistics, or Natural Language Processing (NLP), is not a new field.<sup>7</sup>, sin embargo no es sencillo definir los límites de la disciplina. Así que podemos considerarla como un conjunto de problemas relacionados con fenómenos lingüísticos y una amalgama de soluciones computacionales, de distinto tipo dependiendo del origen del investi-

Según xkcd, 8 los lingüistas computacionales han vivido muy bien hasta ahora vendiendo motos, así que metámonos con ellos.9

AND THE DUMBEST THING ABOUT EMO KIDS ISTHAT ... I ... YOU KNOW, I'M SICK OF EASY TARGETS. ANYONE CAN MAKE FUN OF EMO KIDS.

YOU KNOW WHO'S HAD IT TOO EASY? COMPUTATIONAL LINGUISTS.



"OOH, LOOK AT ME! MY FIELD IS SO ILL-DEFINED I CAN SUBSCRIBE TO ANY OF DOZENS OF CONTRADICTORY MODELS AND STILL BE TAKEN SERIOUSLY!"



<sup>7</sup> http://www.aclweb.org/aclwiki/index.php ?title=Frequently\_asked\_questions \_about\_Computational\_Linguistics

<sup>8</sup> http://www.xkcd.org/114/

<sup>9</sup> http://www.explainxkcd.com/wiki/index.php/ 114:\_Computational\_Linguists