# Advanced Methods for Corpus Processing

Lluís Padró <padro@lsi.upc.edu>
Lluís Màrquez <lluismv@amazon.com>
German Rigau <german.rigau@ehu.es>

Centre de Tecnologies i Aplicacions
del Llenguatge i la Parla

# Foreword

- UPC: Artificial Inteligence
    - EMNLP: Empirical Methods for NLP
    - 2001/2002-2002/2003
        - Lluís Padró, Lluís Màrquez, German Rigau
    - 2003/2004-2004/2005
        - Lluís Padró, Lluís Màrquez, Neus Català, German Rigau
    - 2005/2006
        - L. Padró, L. Màrquez, X. Farreres, J. Daudé, G. Rigau

- EHU: NLP
    - Advanced Methods for Corpus Management
    - 2004/2005-...-2019/2020
        - Lluís Padró, Lluís Màrquez, German Rigau

# Content

- Theme 1: Introduction to corpus analysis.
- Theme 2: Knowledge-based methods.
- Theme 3: Statistical methods.
- Theme 4: Machine learning methods.

# Content

- Knowledge Based methods for NLP (German Rigau)
  - 27 May: 15:00h – 18:00h
  - 28 May: 15:00h – 17:30h
  - 29 May: 15:00h – 17:30h
  - 01 June: 15:00h – 17:30h

- Statistical methods for NLP (Lluís Padró)
  - 2-4 June: 15:00h – 19:00h

- Machine Learning for NLP (Lluís Màrquez)
  - 8-10 June: 15:00h – 19:00h

- Presentations and concluding remarks (German Rigau)
  - ??: 15:00h – 18:00h

# Content

- Knowledge-based NLP (German Rigau)

  - Words & Works
  - Large-scale Knowledge Bases:
    - WordNet & EuroWordNet
  - More large-scale resources
    - ConceptNet, Framenet, VerbNet, PropBank, Predicate Matrix
  - WordNet extensions:
    - SUMO ontology, eXtended WordNet, MCR
  - Ontologies:
    - AdimenSUMO
    - Reasoning, abduction

- Concluding remarks (German Rigau)
  - Combining approaches

# Content

- Statistical methods for NLP (Lluís Padró)

    - Introduction (statistical vs non-statistical NLP, what are statistical models, what is model estimation)
    - MLE Estimation
    - MaxEnt Estimation
    - Hidden Markov Models
    - Structured prediction (sequences): Log-linear models, MEMM, CRF, Perceptron
    - Generalizing structured prediction (dependency structures)

# Content

- Machine Learning for NLP (Lluís Màrquez)

    - Introduction:  Machine Learning and Machine Learning for NLP
    - Machine Learning: Classical Methods from AI
    - Margin-based Machine Learning Algorithms
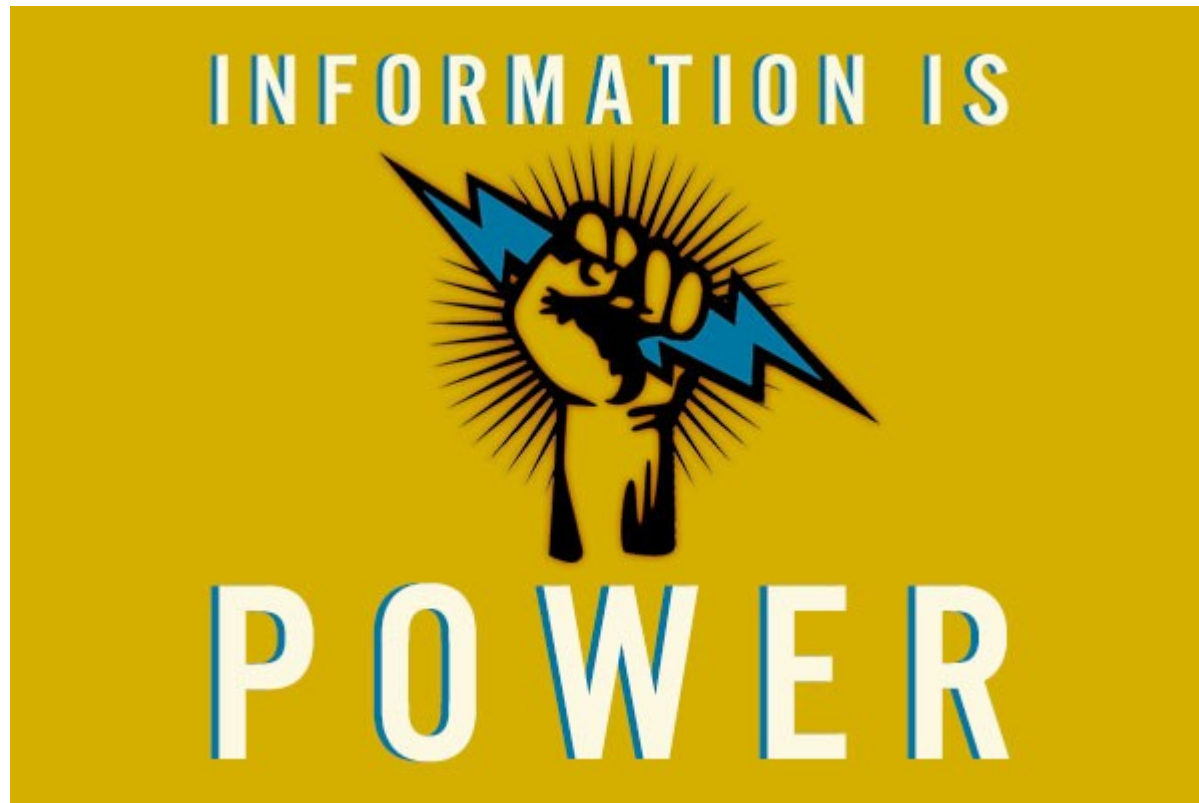    - Machine Learning for NLP
    - Applications

# Evaluation

- Application of empirical methods for NLP:
  - Teacher exercises (30%)

  - Teacher/student topic
  - Short presentation (20%)
    - 15 minutes sharp, ~ 10 slides
    - Presentation: 20/06
  - Written report (50%)
    - Format: http://www.acl2019.org
    - Deadline Report: 20/07
    - Short paper describing an <u>experimental work</u>
      - ~ 3000 words

# Short Motivation

**<span style="color:darkred">Information</span> is power!**

# Short Motivation

# Short Motivation

**Knowledge** is power!

# Short Motivation

# Short Motivation

**<span style="color:#8B0000">Knowledge</span> is power!**

… and the knowledge to use …

# Foreword

*"Cuando creíamos que teníamos todas las respuestas,*
*de pronto, cambiaron todas las preguntas."*

- Mario Benedetti

*"When we thought we had all the answers,*
*suddenly, they changed all the questions. "*

- Mario Benedetti

# Foreword

- Where are the **answers** to the new (and old) questions?

  - Introspection? Experts?…
  - From many people? … "Wisdom of the Crowd"?
  - Books, News, Tweets, … Textual Sources?
  - Multimedia sources? Images, Radio, TV ...
  - Sensors? IoT? ...
  - Anything? Everything?

- Information **overload** …

# Foreword

- **Information overload** …

  - infobesity, infoxication!

# Foreword

- **Information overload** …

  - infobesity, infoxication!
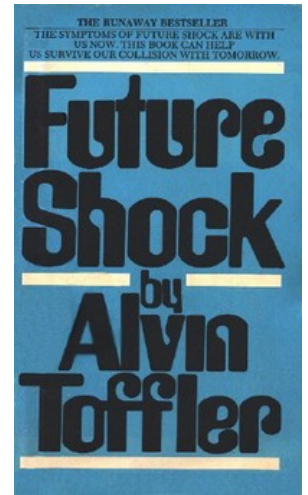  - by Bertram Gross, <u>The Managing of Organizations: The administrative struggle</u> (1964)

# Foreword

- **Information overload** …

  - infobesity, infoxication!
  - by Bertram Gross, *The Managing of Organizations: The administrative struggle* (1964)
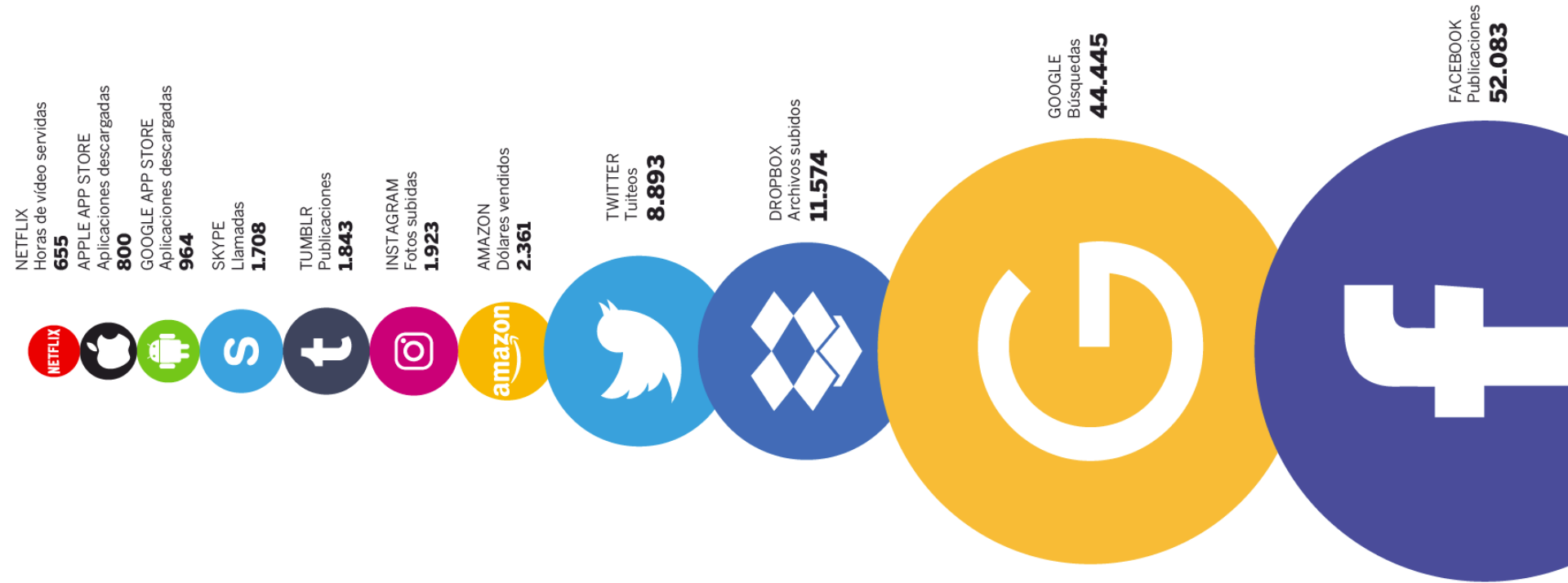  - by Alvin Toffler, *Future Shock* (1970)

# Foreword

- **Information overload** …

  - infobesity, infoxication!
  - by Bertram Gross, _The Managing of Organizations: The administrative struggle_ (1964)
  - by Alvin Toffler, _Future Shock_ (1970)

  - Seneca complained that "_the abundance of books is distraction_" in the 1st century AD!

# Foreword

- **Information overload** occurs when the amount of input to a system exceeds its processing capacity.

- Decision makers have fairly **limited** cognitive processing capacity.

- Consequently, when information overload occurs, it is likely that a **reduction** in decision quality will occur.

- From (Speier et al 1999)


- Always when **advances in technology** have increased a production of information.

# Foreword

- What happens in Internet every **second**? (July 2015)

# Foreword

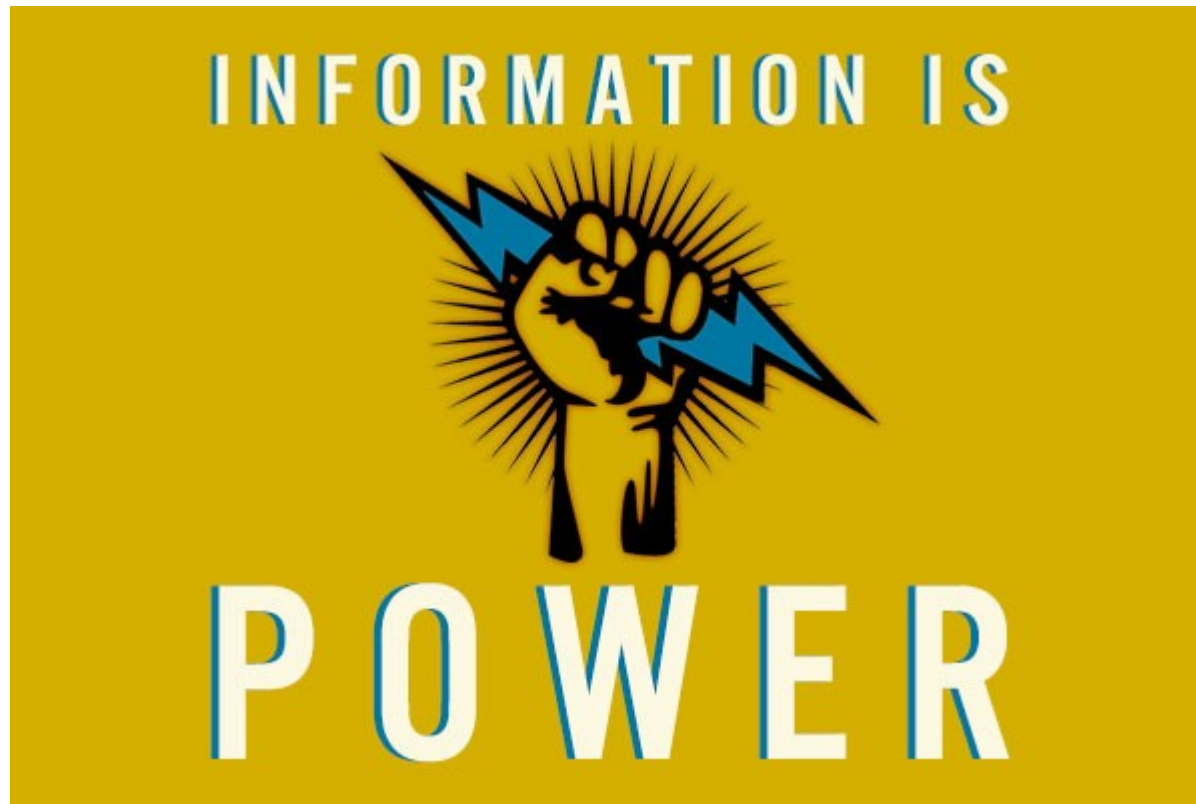- What happens in Internet every **second**? (July 2015)

# Foreword

- … not only coming from Social Media.

- LexisNexis receives **daily** 1.5M news.
- CENDOJ stores 6M judicial sentences (0.6M/year)
- 5M Electronic Health Records (EHR) …
- 0.2M Patents …
- …

- … all kinds of e-documents …

# Foreword

- **Unstructured** digital content accounts for **90%** of all information [White paper IDC 2014] …

- Usually in the form of **texts** and documents in **multiple languages** …

- **Only** appropriate NLP tools can access this wealth of knowledge …

- NLP among the **top** 10 strategic technology trends for 2017 according to Gartner

# Foreword

Because everybody knows that …

# Foreword

But in fact …

# Foreword

e.g. IBM Watson …



but also Google, Facebook, Amazon, Microsoft, …

Advanced Methods for Corpus Processing

# Big Data & NLP ...



http://emm.newsbrief.eu/

# Big Data & NLP …



**http://www.lockheedmartin.com/us/products/W-ICEWS.html**

Advanced Methods for Corpus Processing

# Big Data & NLP …



2012-09-25_0000_US_KNBC_Channel_4_News.jpg

**https://sites.google.com/site/distributedlittleredhen**

# Advanced Methods for Corpus Processing

Lluís Padró <padro@lsi.upc.edu>
Lluís Màrquez <lluismv@amazon.com>
German Rigau <german.rigau@ehu.es>