

CORSO DI  
*WEB MINING E RETRIEVAL*  
*- INTRODUZIONE AL WM -*

---

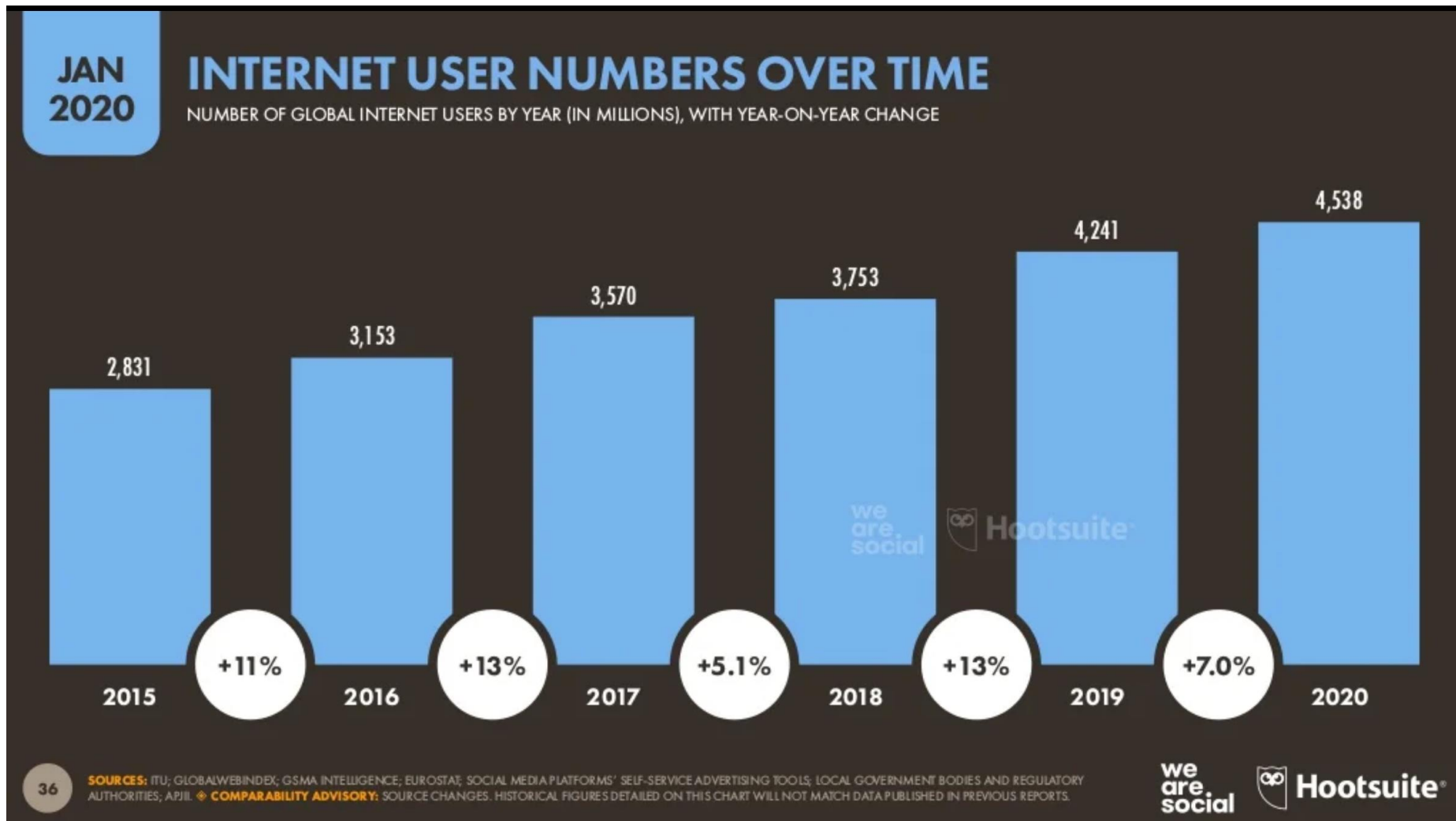
Corso di Laurea in Informatica, Ing. Internet,  
Ing. Informatica, Ing. Gestionale  
(a.a. 2021-2022)

Roberto Basili

# Overview

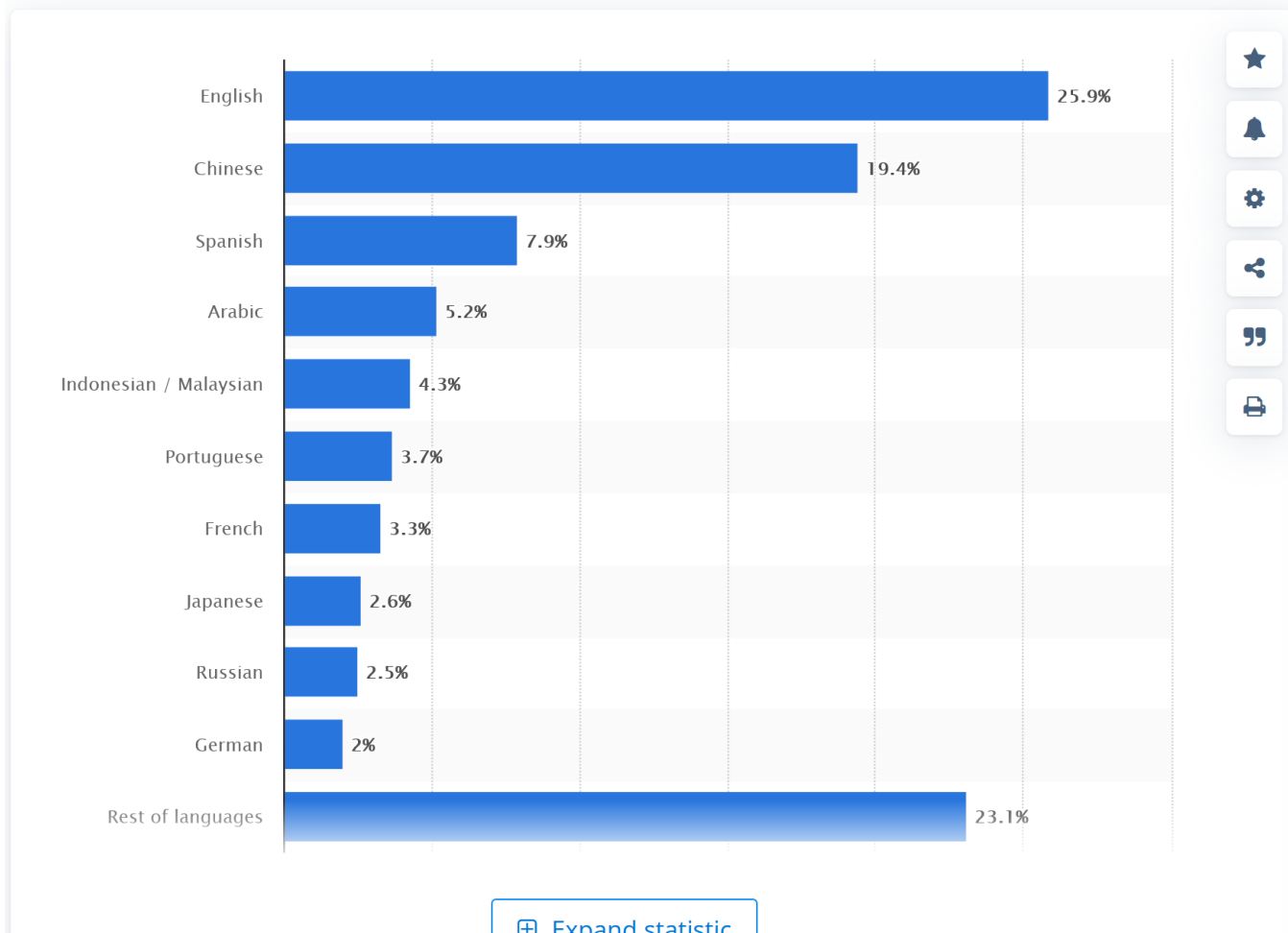
- Web Mining & Retrieval: Motivations & perspectives
  - Web, User-generated contents, Social Media
  - The role of *learning*
  - What is Machine Learning?
  - Data-driven algorithms: sources of complexity
- Main Applications
  - Intelligent Web Search
  - User Profiling for Marketing or Brand reputation management
  - Web Recommending
  - Spoken Dialogue Interaction in Robotics or in Web/mobile Interfaces

# Internet statistics (Jan 2021)

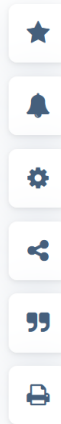
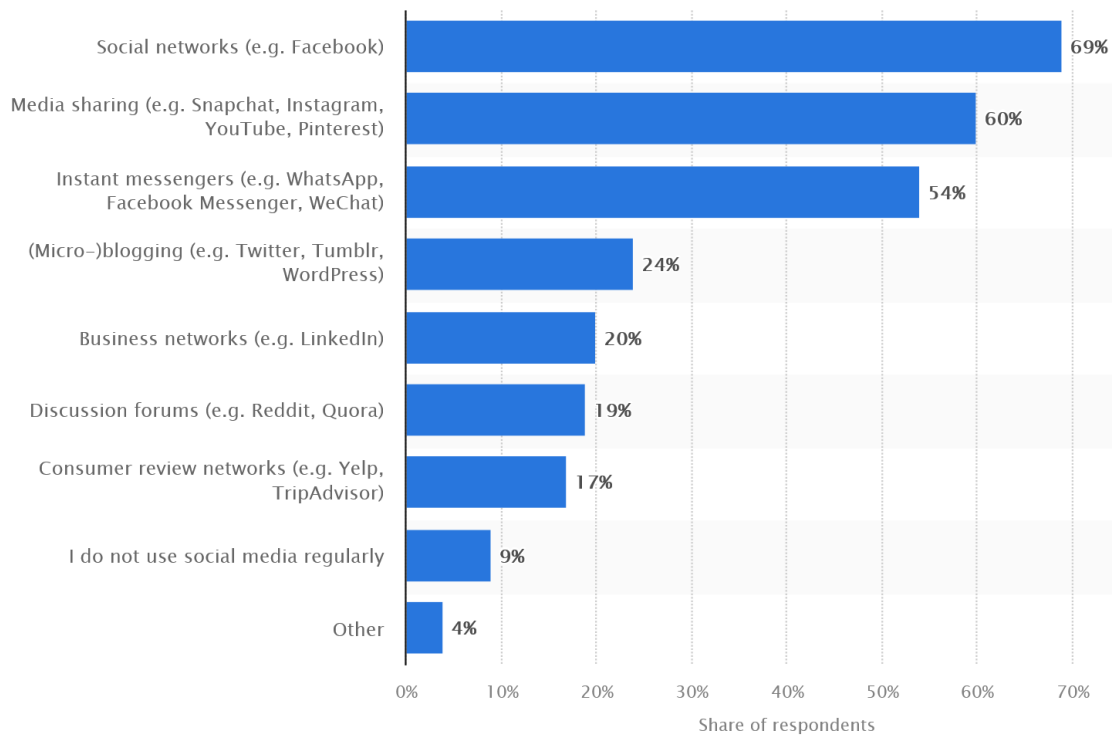


# Internet Statistics (Jan 2020)

Most common languages used on the internet as of January 2020, by share of internet users



## What kinds of social media do you use regularly?



© Statista 2022 🇩🇪

Details:

[Show source](#) ⓘ

Do you know

More than  
**4,000 new books**  
are published every day



Do you know

Contains more  
information than a  
person was likely to  
come across  
**in a lifetime** in the  
**18th century...**



**JAN  
2022**

## ESSENTIAL DIGITAL HEADLINES

OVERVIEW OF THE ADOPTION AND USE OF CONNECTED DEVICES AND SERVICES



GLOBAL OVERVIEW

TOTAL  
POPULATION



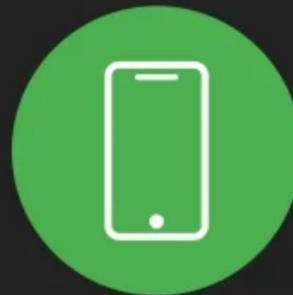
we  
are  
social

**7.91  
BILLION**

URBANISATION

**57.0%**

UNIQUE MOBILE  
PHONE USERS



**5.31  
BILLION**

vs. POPULATION

**67.1%**

INTERNET  
USERS



**4.95  
BILLION**

vs. POPULATION

**62.5%**

ACTIVE SOCIAL  
MEDIA USERS



**4.62  
BILLION**

vs. POPULATION

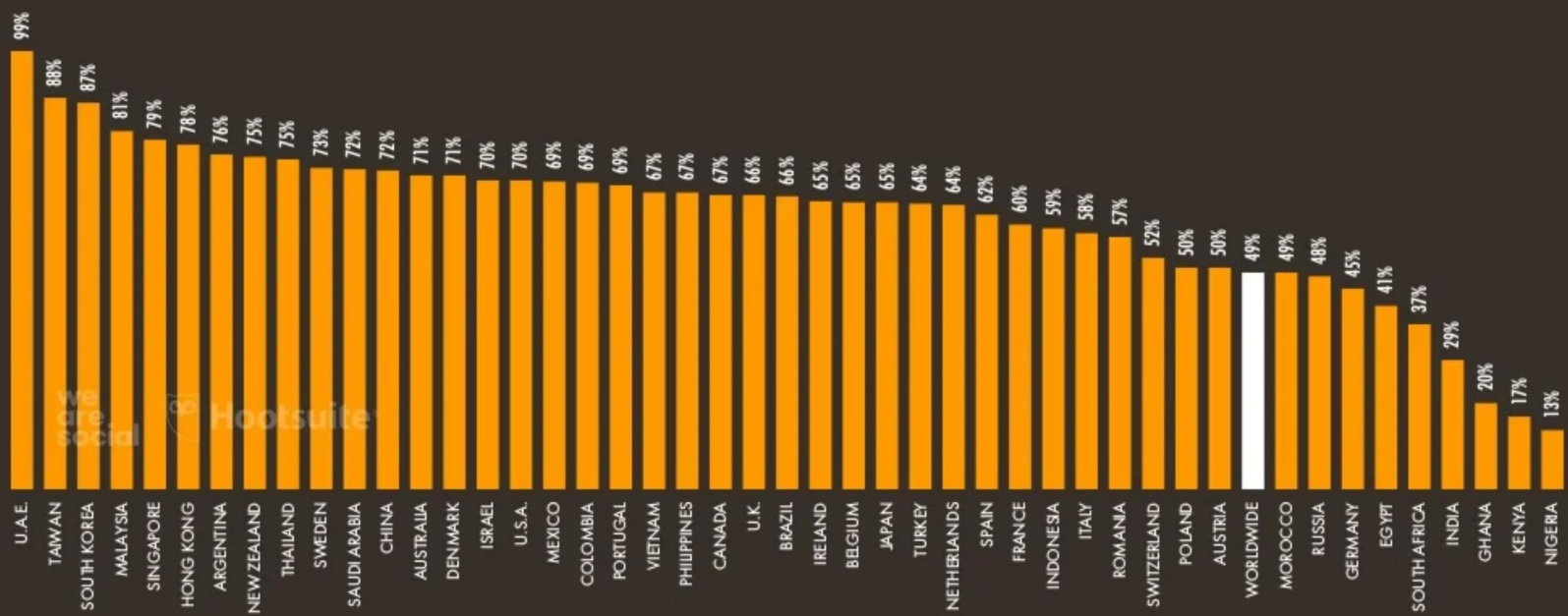
**58.4%**



JAN 2020

# SOCIAL MEDIA PENETRATION

THE NUMBER OF ACTIVE SOCIAL MEDIA USERS COMPARED TO TOTAL POPULATION, REGARDLESS OF AGE



81

SOURCES: KEPIO ANALYSIS; COMPANY STATEMENTS AND EARNINGS ANNOUNCEMENTS; SOCIAL MEDIA PLATFORMS' SELF-SERVICE ADVERTISING TOOLS; MEDIASCOPE; CAFEBAZAR (ALL LATEST DATA AVAILABLE IN JANUARY 2020). \*NOTES: PENETRATION FIGURES ARE FOR TOTAL POPULATION, REGARDLESS OF AGE. ♦ COMPARABILITY ADVISORY: SOURCE AND BASE CHANGES.

JAN  
2020

## DAILY TIME SPENT USING SOCIAL MEDIA

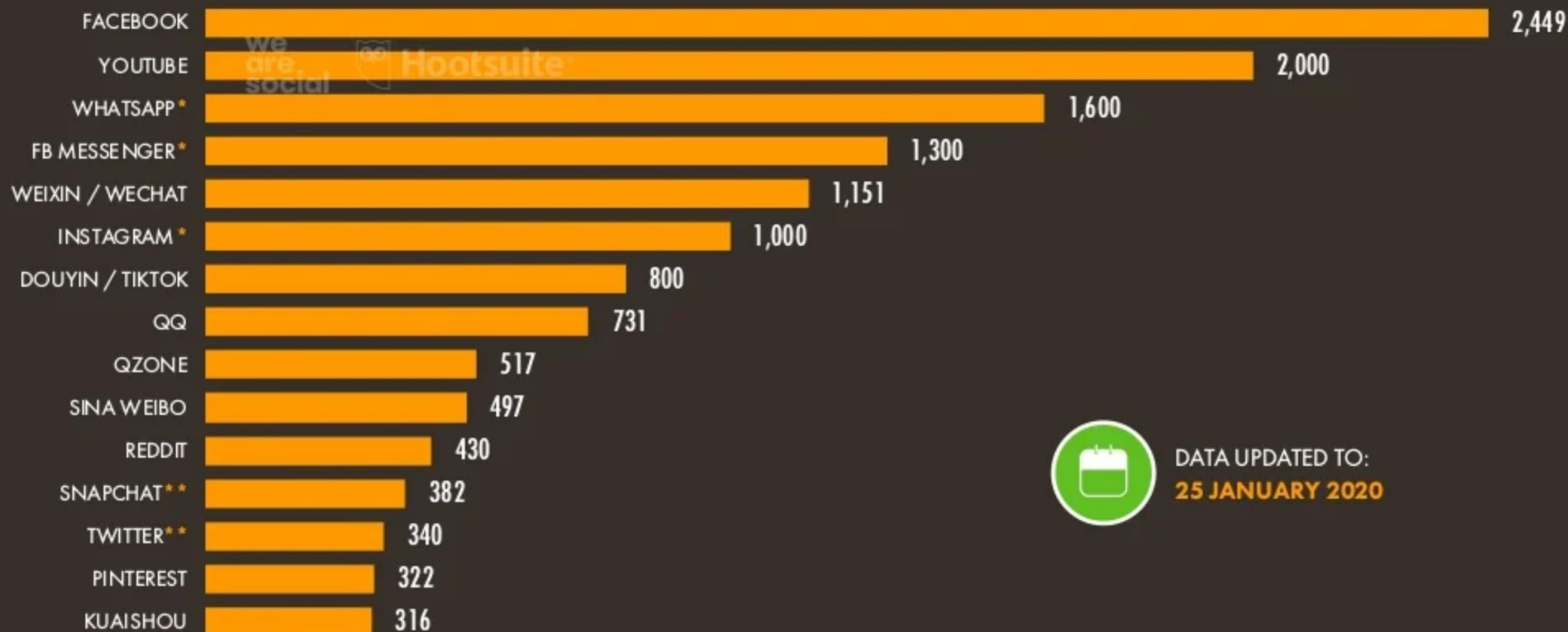
AVERAGE DAILY TIME (IN HOURS AND MINUTES) THAT INTERNET USERS AGED 16 TO 64 SPEND USING SOCIAL MEDIA ON ANY DEVICE



**JAN  
2020**

## THE WORLD'S MOST-USED SOCIAL PLATFORMS

BASED ON MONTHLY ACTIVE USERS, ACTIVE USER ACCOUNTS, ADVERTISING AUDIENCES, OR UNIQUE MONTHLY VISITORS (IN MILLIONS)



DATA UPDATED TO:  
**25 JANUARY 2020**

**JAN  
2022**

# MAIN REASONS

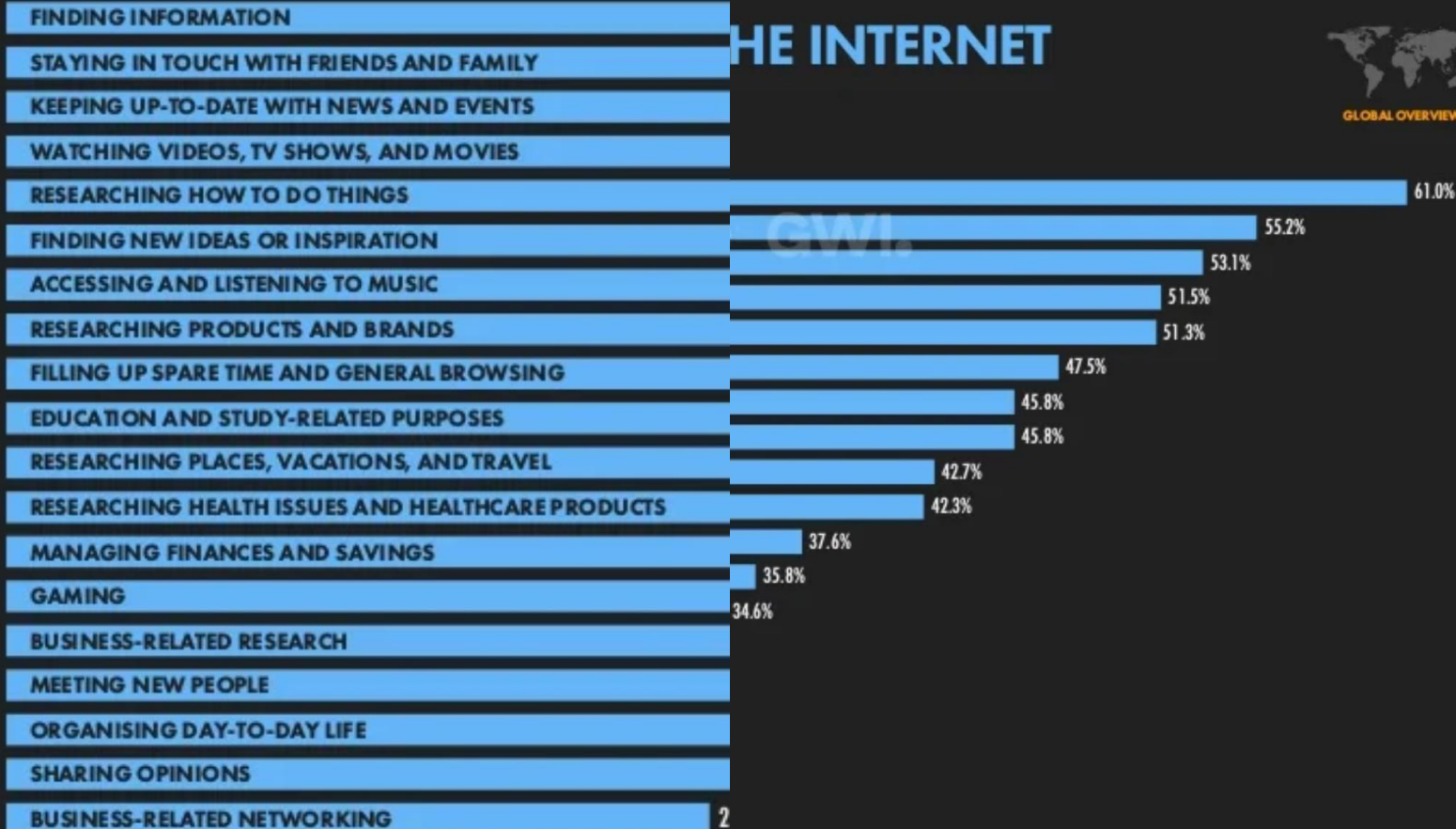
PRIMARY REASONS WHY INTERNET USERS AGED

**JAN  
2022**

## THE INTERNET



- FINDING I
- STAYING I
- KEEPING U
- WATCHIN
- RESEARH
- FINDING I
- ACCESSIN
- RESEARH
- FILLING U
- EDUCATIO
- RESEARH
- RESEARH
- MANAGIN
- GAMING
- BUSINESS
- MEETING I
- ORGANIS
- SHARING
- BUSINESS





JAN  
2020

## TWITTER AUDIENCE OVERVIEW

THE POTENTIAL NUMBER OF PEOPLE THAT MARKETERS CAN REACH USING ADVERTS ON TWITTER

NUMBER OF PEOPLE THAT  
TWITTER REPORTS  
CAN BE REACHED WITH  
ADVERTS ON TWITTER



**339.6**  
MILLION

SHARE OF POPULATION  
AGED 13+ THAT MARKETERS  
CAN REACH WITH  
ADVERTS ON TWITTER



**5.6%**

QUARTER-ON-  
QUARTER CHANGE  
IN TWITTER'S  
ADVERTISING REACH



**-3.1%**

PERCENTAGE OF  
ITS AD AUDIENCE  
THAT TWITTER  
REPORTS IS FEMALE\*



**38%**

PERCENTAGE OF  
ITS AD AUDIENCE  
THAT TWITTER  
REPORTS IS MALE\*



**62%**

**JAN  
2020**

## MOST-USED EMOJI ON TWITTER

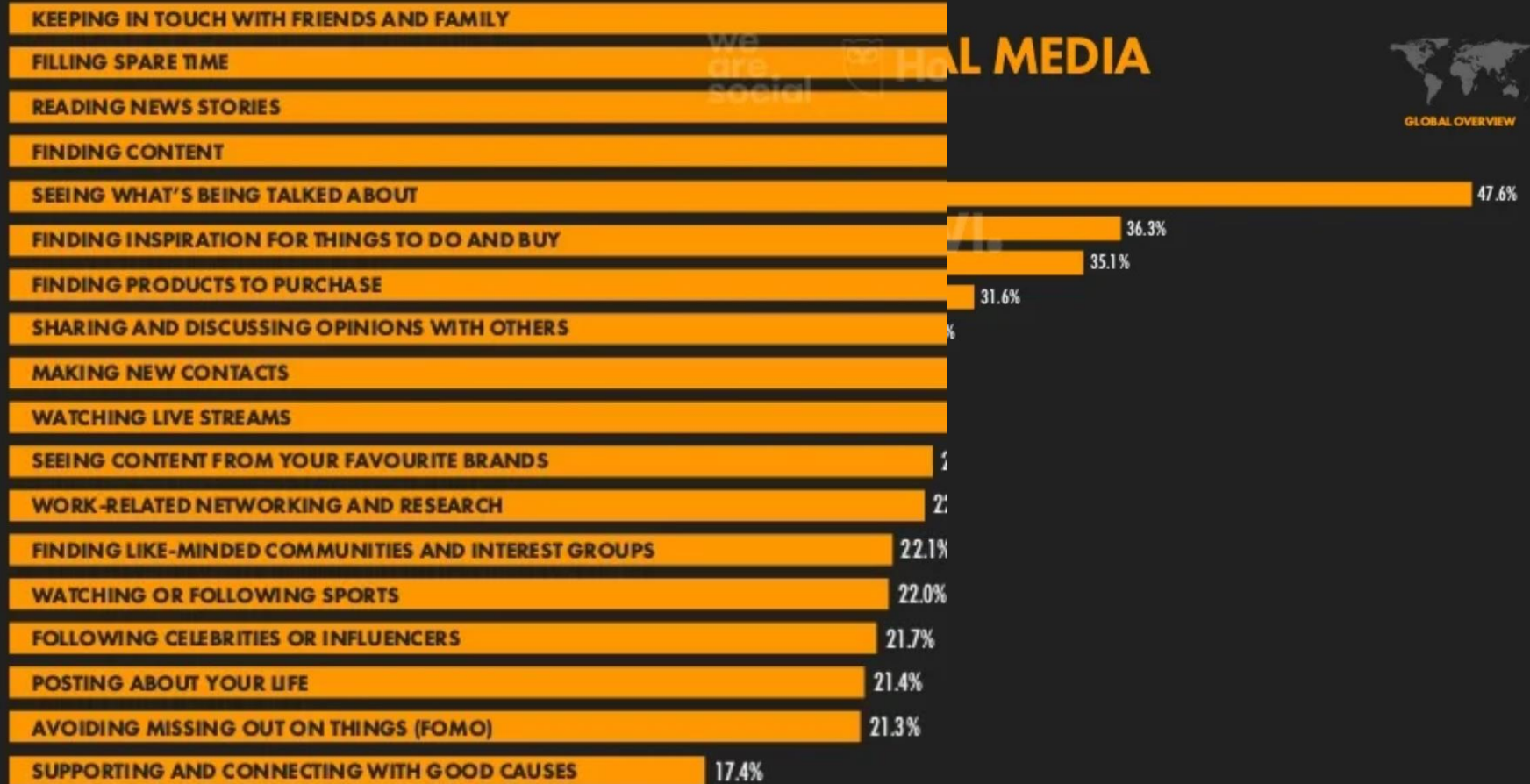
EMOJI THAT HAVE BEEN USED THE GREATEST NUMBER OF TIMES ON TWITTER (ALL TIME)

#	EMOJI	TIMES USED	#	EMOJI	TIMES USED	#	EMOJI	TIMES USED	#	EMOJI	TIMES USED
01		2,671,000,000	11		428,000,000	21		245,000,000	31		198,000,000
02		1,289,000,000	12		389,000,000	22		238,000,000	32		193,000,000
03		966,000,000	13		382,000,000	23		237,000,000	33		191,000,000
04		964,000,000	14		365,000,000	24		236,000,000	34		187,000,000
05		817,000,000	15		359,000,000	25		232,000,000	35		182,000,000
06		743,000,000	16		336,000,000	26		229,000,000	36		181,000,000
07		632,000,000	17		309,000,000	27		217,000,000	37		168,000,000
08		500,000,000	18		273,000,000	28		216,000,000	38		165,000,000
09		493,000,000	19		258,000,000	29		212,000,000	39		163,000,000
10		475,000,000	20		246,000,000	30		199,000,000	40		163,000,000

**JAN  
2022**

# MAIN REASONS FOR USING SOCIAL MEDIA

PRIMARY REASONS WHY INTERNET USERS AGED 16 TO 64 USE SOCIAL MEDIA



we  
are  
social

GLOBAL MEDIA



GLOBAL OVERVIEW

# WE ARE SOCIAL'S PERSPECTIVE: SOCIAL IN 2020

## SHIFTS IN HOW PEOPLE BEHAVE AND INTERACT ON SOCIAL



### BAD INFLUENCE

Being a creator has lost its lo-fi sheen; many lifestyle influencers lead unrelatable lives, while celebrity 'creators' like Will Smith are blowing up on platforms like YouTube and TikTok. As a result, there's a growing backlash against influencer culture and the metrics that drive it.

**In 2020, brands will look beyond likes, followers and reach to generate genuine engagement**



### ADDED VALUE

The internet has long been a wild west where intellectual property is barely there. But in a maturing digital frontier, creators have grown dedicated audiences who not only see value in their content, but recognise their style anywhere. As a result, communities are rallying to protect creators.

**In 2020, brands will take greater steps to ensure they're being respectful of digital communities**



### RUNNING COMMENTARY

Audiences are increasingly willing to invest time and attention in content and narratives they deem to have a higher value. This isn't about a shift back to traditional media. It's about longer, more complex content designed to be consumed in-platform and on smaller screens.

**In 2020, brands will tell more complex stories across multiple touchpoints on social**



# Dealing with *real* Social media data



# WM&R: Motivations

- *What does Web Mining mean?*
- *Why Information Retrieval is involved?*
- *Why Machine Learning?*
- *Which are the contributions of IR/ML to technologies that support and exploit Web Data, Information and Knowledge?*
- *Which are the technological perspectives in the medium-long term?*

# What is Web Mining?

- *Web Mining* refers to a body of technologies currently needed for the *exploitation of publicly available information from the Web and the IoT*
  - Contents: data but also ... people, locations, events, concepts, ...
  - Relations:
    - Links within structured networks
    - Thematic, interpersonal and semantic associations
    - Analogies
  - On-Line Structured and semi-structured resources (e.g. Wikipedia)
  - Textual, Multimedia and Multilingual Contents
  - Trends e time-related information (community on-line behaviours)
  - Opinions, Preferences, Expectations

# Why IR?

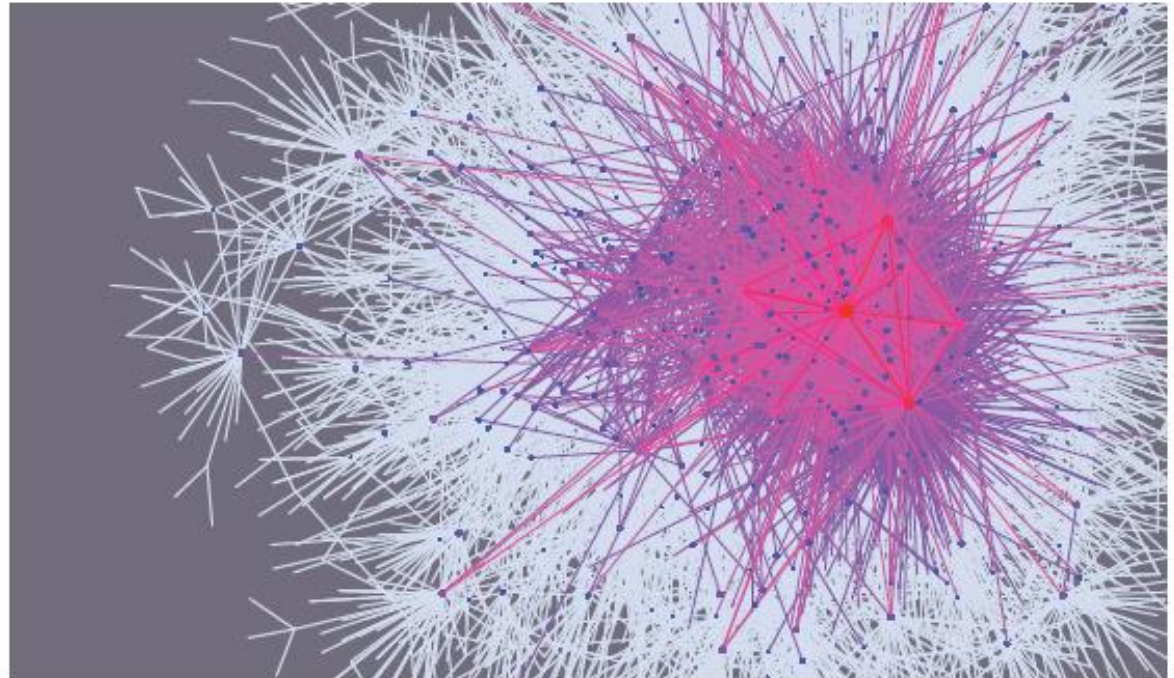
- The volumes involved in Web Mining pose the crucial problem of *locating information* beforehand
- Automatic information access is possible only if we solve the two major challenges
  - **What** is relevant
  - **Where** the relevant information is located
- **Searching information corresponds to computing an uncertain function that models the mapping between information needs and the targeted data**

# Machine Learning vs IR?

- Web mining involve heterogeneous information that is characterized search as strongly uncertain process
- The available information is characterized by:
  - Incompleteness:
    - Short queries as an incomplete description of the information need
  - Variability: Wealth of data vs. heterogeneity of formats and access modes
    - Contents are dispersed in various forms across data sources
  - Vague Requirements
    - Information is often implicit (i.e. partially and qualitatively expressed) in the operational contexts
  - Subjectivity
    - Relevance depends on the user and not just on the contents
  - Timeliness
  - Authority

# Machine Learning vs. IR

- Uncertainty is so pervasive that exhaustive solutions (i.e. global *optima*) are not available or even not existing
- “*Finding diamonds in the rough*”  
(Fan Chung, UCSD)





# Machine Learning vs. IR

- Le tecniche di ML propongono una ampia serie di algoritmi, strategie e tecniche per la produzione di soluzioni *sub-ottime* ma efficaci
- Nel processo di *learning* i dati suggeriscono la ipotesi risolutiva per la funzione di *mapping*
- Tale ipotesi è attesa migliorare la prestazione complessiva del sistema di base
  - Accuratezza
  - Efficienza computazionale

# Machine Learning

- (Langley, 2000): l'Apprendimento Automatico si occupa dei meccanismi attraverso i quali un agente intelligente migliora nel tempo le sue prestazioni  $P$  nell'effettuare un compito  $C$ .
- La prova del successo dell'apprendimento è quindi nella capacità di misurare l'incremento  $\Delta P$  delle prestazioni sulla base delle esperienze  $E$  che l'agente è in grado di raccogliere durante il suo ciclo di vita.
- La natura dell'apprendimento è quindi tutta nella caratterizzazione delle nozioni qui primitive di *compito*, *prestazione* ed *esperienza*.



# Esperienza ed Apprendimento

- L'esperienza, per esempio, nel gioco degli scacchi può essere interpretata in diversi modi:
  - i dati sulle vittorie (e sconfitte) pregresse per valutare la bontà (o la inadeguatezza) di strategie e mosse eseguite rispetto all'avversario.
  - valutazione fornita sulle mosse da un docente esterno (oracolo, guida).
  - Adeguatezza dei comportamenti derivata dalla auto-osservazione, cioè dalla capacità di analizzare partite dell'agente contro se stesso secondo un modello esplicito del processo (partita) e della sua evoluzione (comportamento, vantaggi, ...).

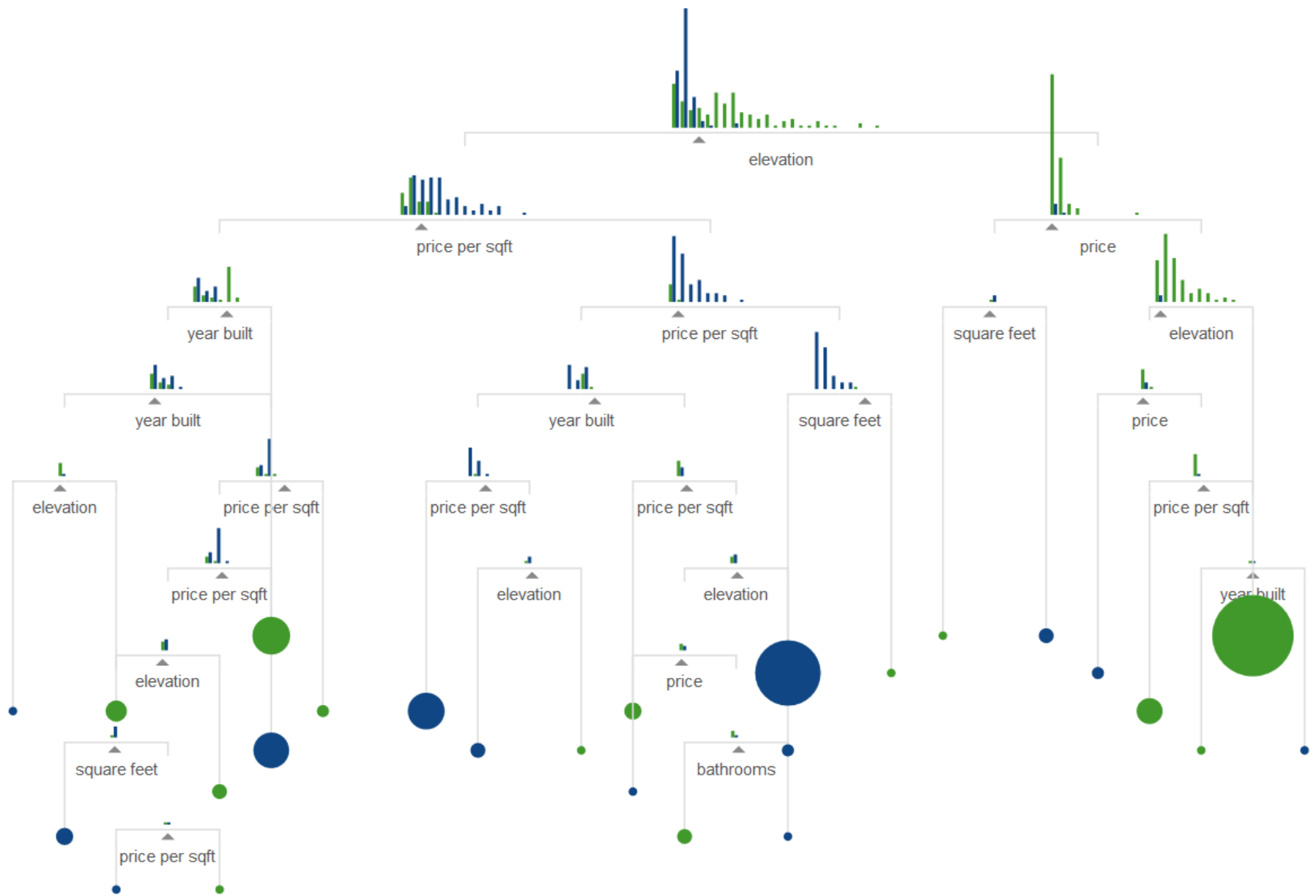
# ML: una introduzione visuale

- See URL: [http://www.r2d3.us/visual-intro-to-machine-learning-part-1/?imm\\_mid=0d76b4&cmp=em-data-na-na-newsltr\\_20150826](http://www.r2d3.us/visual-intro-to-machine-learning-part-1/?imm_mid=0d76b4&cmp=em-data-na-na-newsltr_20150826)

# Apprendimento e Classi di Algoritmi

- Acquisizione di:
  - Funzioni logiche booleane, (ad es., alberi di decisione)
  - Induzione: determinazione ricorsiva delle CNES che caratterizzano i diversi sottogruppi .
- Approcci probabilistici:
  - Funzione target di Probabilità, (ad es., classificatore Bayesiano)
  - Induzione: Stima delle probabilità (in quanto parametri).
- Approcci geometrici
  - Funzioni di separazione in spazi vettoriali (lineari e non)
    - KNN
    - Funzioni Lineari, perceptroni, Neural Networks, Support Vector Machines,...
    - Embeddings, analisi spettrale (trasformazioni di spazio)
  - Induzione: parametrizzare la funzione appartenente ad una certa classe (ad es. polinomi di grado  $n$ )

# Es. apprendimento alberi di decisione

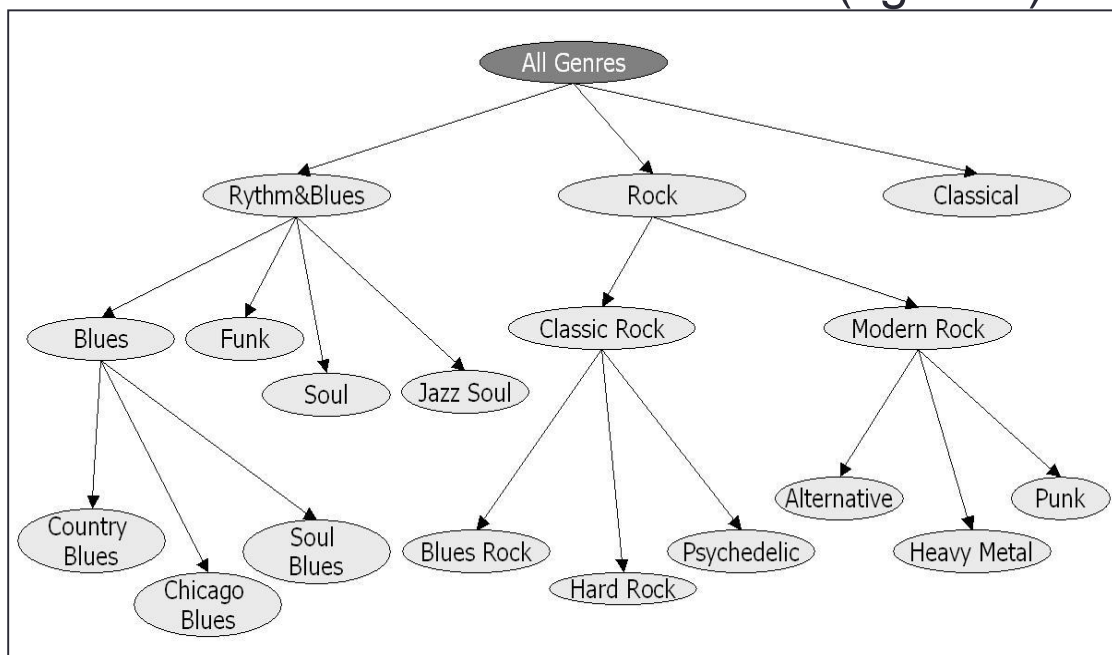


# Apprendimento senza supervisione

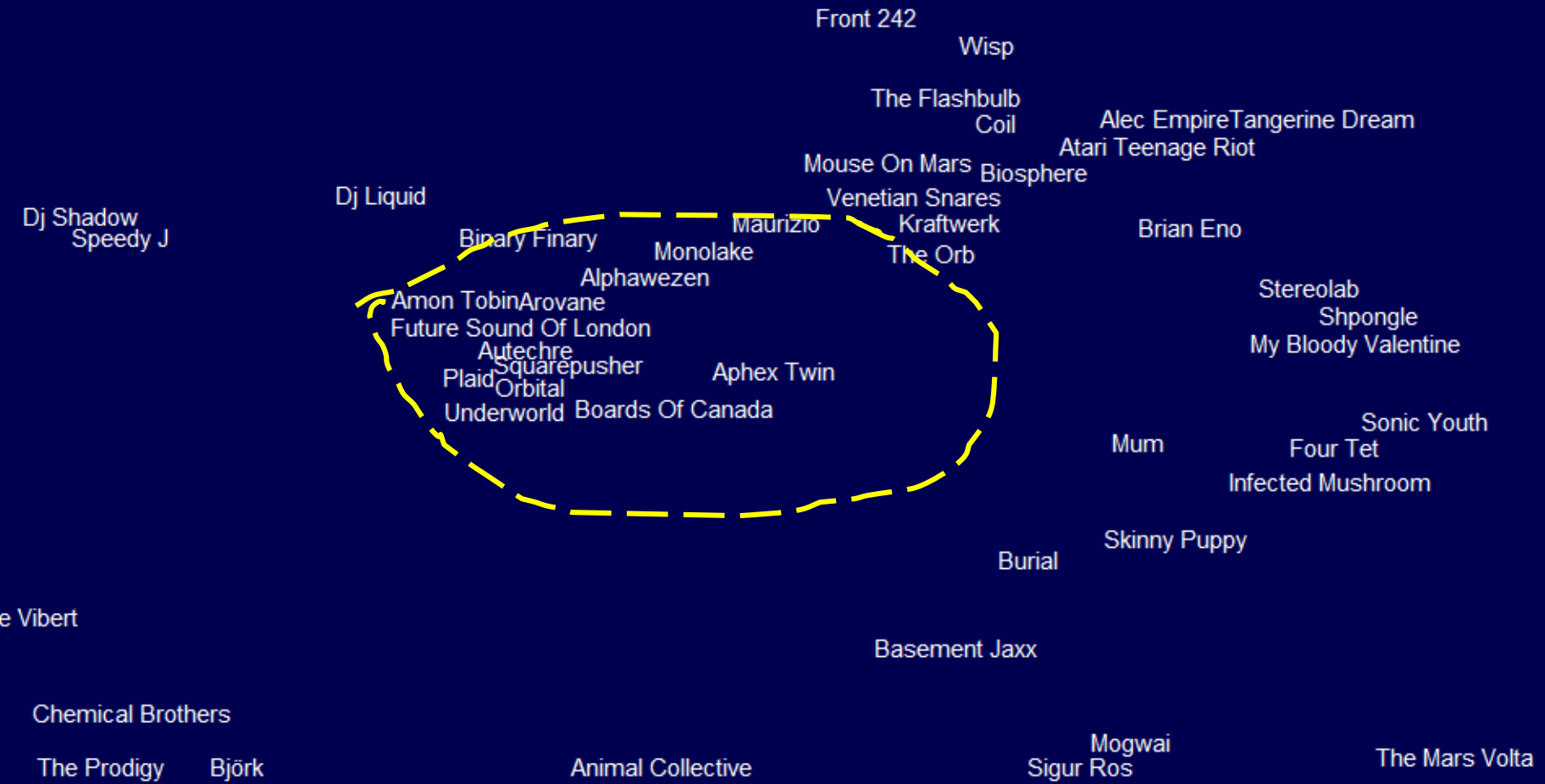
- In assenza di un oracolo o di conoscenze sul task esistono ancora molti modi di migliorare le proprie prestazioni, ad es.
  - Migliorando il proprio modello del mondo (acquisizione/*discovery* della conoscenza)
  - Migliorando le proprie prestazioni computazionali (ottimizzazione)

# Apprendimento senza supervisione

- Esempio: Al termine del processo di acquisizione il sistema dispone di un sistema di classi e relazioni indotti che migliora la sua interazione futura con l'ambiente operativo (ad es. l'utente)
- Il miglioramento avviene quindi almeno rispetto agli algoritmi di ricerca: la organizzazione gerarchica consente di esaminare solo i membri dell'insieme in alcune classi (i generi).

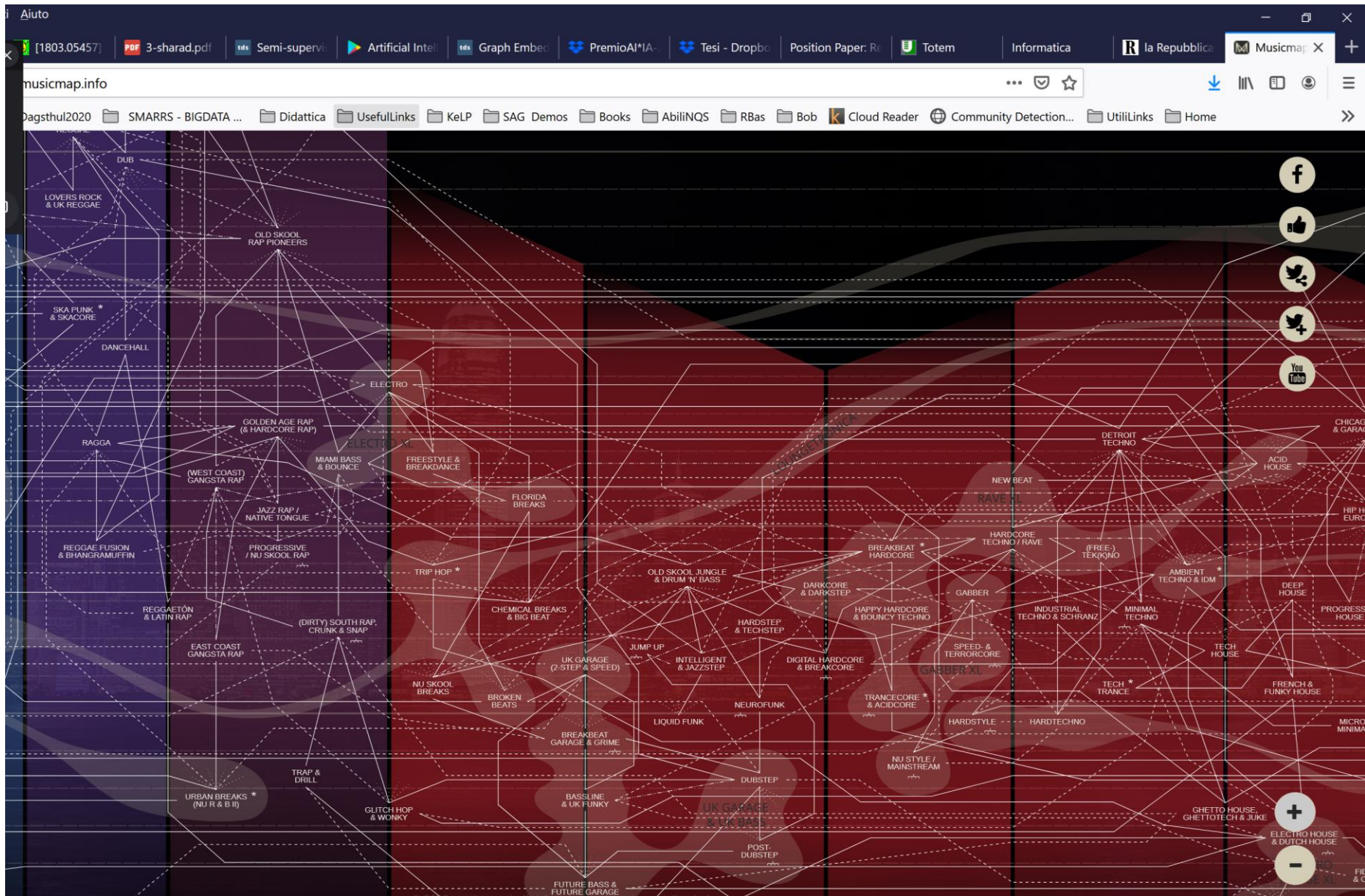


# map





# Music maps: 2020





# Web IR

- Processing Web data: content detection, link detection, ...
- Web Crawling
- Web Search: indici, link analysis
- Ranking: weighting contents, links and formats, authority, timeliness
- Meta-search
- Link Analysis

# Information, Web and Natural Language

Web contents, characterized by rich multimedia information, are mostly **opaque from a semantic standpoint**

Firefox - Mozilla Firefox

File Modifica Visualizza Cronologia Segnalibri Yahoo! Strumenti Aiuto

Errore caricamento pagina blunauta baloon via casilina neg... AGI China 24 - Quotidiani on line Errore caricamento pagina

Indietro Avanti DownloadHelper www.takungpao.com.hk

Più visitati Music RAI Meteo Italia - Previsioni ... Rivista Universitas La Repubblica.it - Hom...

QuickStores Cerca Tutti Su eBay Su Amazon Su Shopping.com

blunauta baloon Via Casilina ne SEARCH 49°

今天是 2011年11月13日 星期日 顯示器最佳分辨率 1024X768 今日天氣 加入最愛 設為首頁 大公網新版

2011 中国証券金紫荆獎 Golden Bauhinia Awards

首頁 國內 國際 港澳 兩岸 評論 財經 體育 教育 科技 醫學 娛樂 文化 副刊 軍事 生活 旅遊 圖片 博客

關鍵詞: 欄目: 全部 最近三個月 三個月之前 檢索 手機新聞 手機博客 漢語學習 新聞點擊排行

滾動新聞:

**胡總語特首:防範經濟金融風險**  
胡錦濤在夏威夷會見出席APEC峰會的曾蔭權。他祝賀香港區議會選舉成功,並充分肯定曾蔭權及港府工作,要求做好經濟金融風險防範

**胡連會登場 共同宣示九二共識**  
胡錦濤第四次在APEC峰會期間會見連戰。他強調,認同「九二共識」是兩岸開展對話協商的必要前提,也是兩岸關係和平發展的重要基礎

西藏黨代會高調反「藏獨」 德國作家:外埠雜誌報道西藏  
傳媒入日本福島核電站探險 英國大裁軍 傷兵難雜免  
滇礦難已30死 13人生還 礦工講述內幕 事故並不意外  
范徐麗泰認民望跌最不耐 選委再獲60提名表 累積逾千人  
聖保羅中學本月底截止招 選委再獲60提名表 累積逾千人  
民調逆轉 藍高層:國親吵鬧地 秋門訴求多 向藍綠表不滿  
世界新七奇觀 亞洲景佔四 新奇觀選舉爭議

中國實體書店苦苦掙扎求 加入TPP 台密集會談探險  
香港人家/蔡仕榮 人生導師 活出自我 我香港人家/教導子女...  
債務危機好 港ADR幾全線造 歐元反降 兌美元逼近1.38  
入世十年/充分對接 華強北最 入世十年/挑戰「二次...  
抽除「雜車」 工人險生 南亞漢命案 警日籍妻

即時新聞

- 組國/河南全國太極拳錦標賽賽
- 奧巴馬重申美不支持「台灣獨立
- 巴基斯坦西北部兩起襲擊 16人
- 圖文/胡錦濤會見美國總統奧巴
- 兩岸30對愛侶在廈門集體證婚
- 中日韓衛生部長會議在青島舉行
- 面向中國遊客中英雜誌紐約創刊
- 「CEO聖經」成內地官員考試
- 斯特恩:經紀人是勞資談判的障
- 香港冀成爲人幣國際化關鍵角色
- 日學者提出地核物質形態新假設
- 中國影視機構向國際大師「取經

焦點關注

區議會選舉 香港

2011APEC 港黑金事件 201

神八天宮對接 第七次陳江會 李

9.1衝擊事件 中國航母試航 辛

http://www.takungpao.com.hk/news/11/11/13/2011\_apec\_xgbd-1423309.htm

# Information, Web and language

Hu meets KMT honorary chairman in Hawaii - People's Daily Online - Mozilla Firefox

File Modifica Visualizza Cronologia Segnalibri Yahoo! Strumenti Aiuto

Hu meets KMT honorary chairman

Indietro Avanti Down

*Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011.*  
*(Xinhua/Huang Jingwen)*

HONOLULU, United States, Nov. 11 (Xinhua) -- Hu Jintao, general secretary of the Central

Latest News: • Indonesia to host European Higher Education Fair

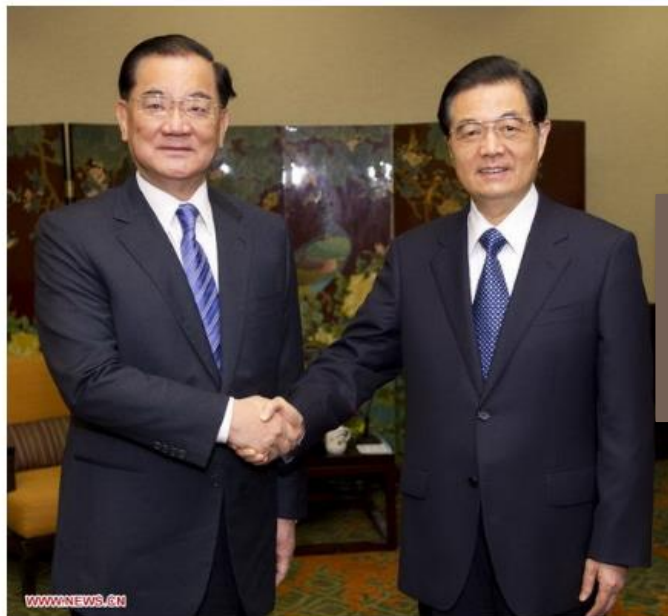
Beijing Sunny 15 / 1 City Forecast

Home >> China Politics

## Hu meets KMT honorary chairman in Hawaii

(Xinhua)

11:10, November 12, 2011 🔍 +-



*Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011.*

Selections for you



Miao ethnic group celebrates Miao's New Year in SW China



World's first Angry Birds exclusive shop opens in Helsinki

## Who is Hu Jintao?

Most Popular

- 1 Hu reaffirms support to Hong Kong's sta...
- 2 Hu meets KMT honorary chairman in Hawaii
- 3 China in APEC: a mutually beneficial en...
- 4 Night life in Shanghai
- 5 China's 2011 foreign trade to grow 20 p...
- 6 Beijing house prices stumble 5.1 pct as...
- 7 Lama students start school in Tibet Col...
- 8 Police in central China crack phoney ca...



Hu Jintao



Ricerca

Circa 725.000 risultati (0,09 secondi)

- Tutto
- immagini**
- Mappe
- Video
- Notizie
- Shopping
- PIÙ conte

Tutti i ri  
Per argomento

- Qualsiasi dimensione
- Grandi
  - Medie
  - icone
  - Maggiori di...
  - Dimensioni esatte...

- Qualsiasi colore
- A colori
  - Bianco e nero
- 

- Qualsiasi tipo
- Volti
  - Foto
  - Clip art
  - Disegni

Visual standard  
Mostra dimensioni



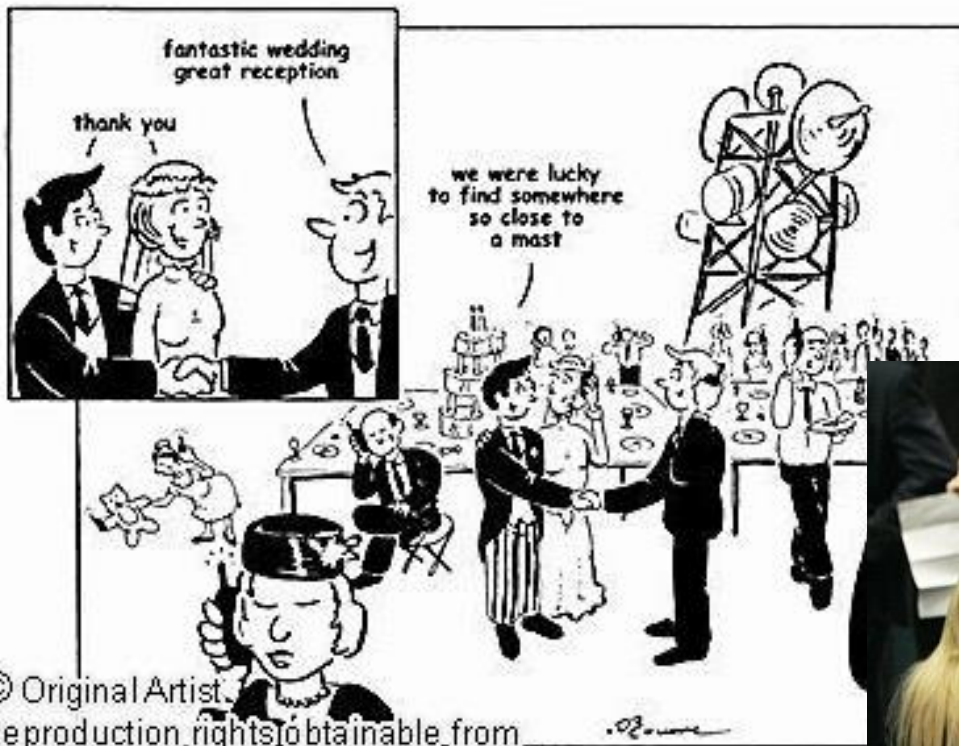


# Content Semantics and Natural Language

- Human languages are the main carrier of the information involved in processes such as *retrieval*, *publication* and *exchange* of knowledge as it is associated to the open Web contents
- Words and NL syntactic structures express concepts, activities, events, abstractions and conceptual relations we usually share through data
- “*Language is parasitic to knowledge representation languages but the viceversa is not true*” (Wilks, 2001)

# NL and Knowledge

- Natural Languages are even too successful in modern



search |

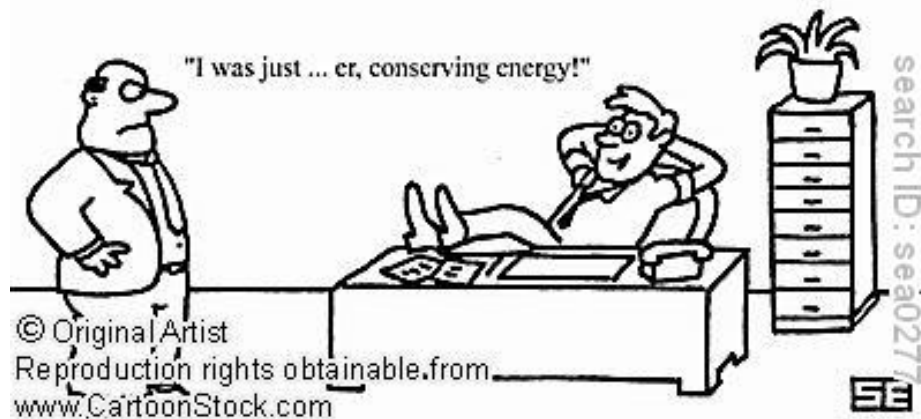


© Original Artist.  
Reproduction rights obtainable from  
[www.CartoonStock.com](http://www.CartoonStock.com)

© AFP/Getty Images

# Benefits of a data-driven approach

- Very effective learning algorithms available (e.g. Support Vector Machines)
- The ML technology is *portable* while imperative coding is *task (i.e. scenario) specific*
- Very accurate solutions can be obtained
- Gathering training data *much less expensive* than rule coding
- In dynamically evolving scenarios, incremental refinement of the system only consists in re-training



# Data Mining: perspectives and benefits

- Technical advantages
  - Self adaptivity to changing operational conditions (i.e. domains)
  - Better SW management and incremental maintenance
  - More flexibility for special-purpose versioning:
    - No need for re-engineering or independent software developments
    - Just new domain-specific examples are needed
- Cost benefits
  - The data-driven approach has been shown to reduce the development costs up to 80-90% in several NLP tasks
- Market benefits
  - Reduced time-to-market
  - Competitive advantages: the lack of similar products makes the system targeted strongly competitive solutions



# Semantics and News

Applicazioni Risorse Sistema mar 27 lug, 23.47 dan

Gmail ... x SRL\_EN x Come ... x R Econo... x Googl... x Tanl It... x Frame... x SRL\_EN x Econo... x

file:///home/danilo/Downloads/SRL\_ITA/sorgente/Economia%20-%20Repubblica.it.html

Telefilm in stream... Flash Forward pri... Cronologia Altri Pr



L'ad punta a nuove regole sulla base del modello Pomigliano. L'annuncio, che prevede l'uscita da Federmeccanica, domani al vertice con il governo o giovedì con una lettera a Bombassei. Potrebbe avvenire assieme alla decisione di creare una new company per

Pomigliano di SALVATORE TROPEA

**Cisl-Uil: "L'accordo di categoria non si tocca"** di S. PAROLA

**Saconi: "Su Fiat partita aperta"**

**Nasce Fabbrica Italia Pomigliano**

## Si dimette il capo di Bp buonuscita un milione di sterline



Oggi l'annuncio: a Tony Hayward subentrerà il direttore esecutivo Robert Dudley. **I costi legati al disastro sono saliti a 32,2 miliardi di dollari**, ma la società li deterrà evitando di versare al fisco Usa 10 miliardi

## Manager Usa, è Ellison di Oracle il più pagato del decennio



Ha guadagnato 1,84 miliardi di dollari. Nella classifica del *Wall Street Journal* sui leader delle società quotate, secondo con 1,14 miliardi il capo di Expedia, terzo Irani di Occidental Petroleum. Solo quarto Steve Jobs

### Il nemico alle porte

La Consob e la mano invisibile

Altri articoli



**PICCOLE GRANDI IMPRESE**  
DI LUCA PAGNI

La grande sfida del teleshopping

La crisi colpisce anche i porti turistici  
ma siamo sicuri che sia un male?

Altri articoli



**PERCENTUALMENTE**  
DI ROSARIA AMATO

La prova del 9

L'export risolve il Pil, ma non le famiglie

Altri articoli

### GLI ESPERTI RISPONDONO

CASA

A cura di Antonella Donati



**Compenso extra, quando ne ha diritto l'amministratore**

Mia moglie ed il fratello sono proprietari di un appartamento in condominio. Allo stato

Il tuo libro arriva dove  
hai sempre sognato.

ilmiolibro.it

**24ORE AGI**

**Roma 19:04**  
ACEA: NEL I SEMESTRE UTILE NETTO +52,1% A 82, MLN

**Parigi 18:42**  
AIR FRANCE-KLM: TORNA IN UTILE NEL PRIMO TRIMESTRE

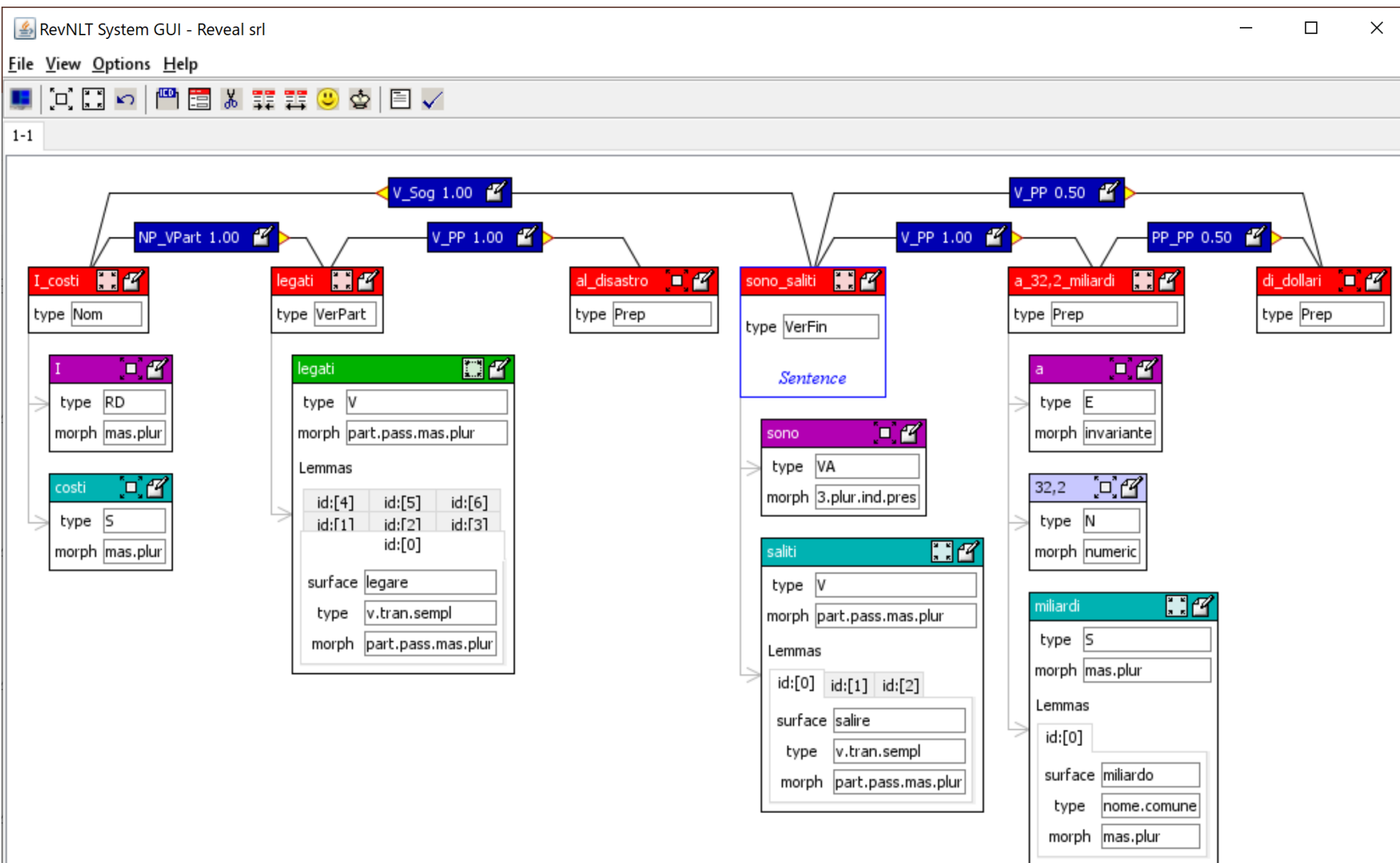
← 3 → Le altre not

**CREDITO ALLE IMPRESE**

**Microimprese: con la crisi aumenta il rischio di credito**

IN COLLABORAZIONE CON

# Ouput del parser NLP



# Laboratori del Corso

- Nell'ambito dei Laboratori agli studenti saranno resi disponibili:
  - Piattaforme di Machine Learning: Weka, KELP, Pytorch, SciKit
  - Motori di Ricerca: Lucene, Solr
  - Strumenti di AI per l'elaborazione dei testi:
    - Recursive Neural Networks per l'apprendimento di lessici vettoriali
    - Transformers per il semantic parsing e la textual inference
    - Parser grammaticali delle lingue (ita, eng)
    - Named Entity recognition and Wikification tools

# Un esempio: Kelp: Java-based kernel framework

**GitHub**

Explore Features Enterprise Pricing

Sign up

Sign in

## KeLP

KeLP (Kernel-based Learning Platform) is a Java machine learning platform developed within the SAG group and the QCRI.

University of Roma, Tor Vergata <http://sag.art.uniroma2.it/demo-software/kelp>

Repositories People 3

Filters Find a repository...

### kelp-additional-algorithms

Updated 9 days ago

### kelp-full

Updated 12 days ago

### Semantic Analytics Group @ Uniroma2

SAG is the Semantic Analytics Group at the University of Rome, Tor Vergata

People Research Teaching Publications Projects Demo & Software Contacts

### KeLP (Kernel-based Learning Platform)

KeLP (Kernel-based Learning Platform) is a machine learning platform developed within the SAG group. It is entirely written in Java and it is strongly focused on *Kernel Machines*. It includes different Online and Batch Learning and Classification algorithms, Kernel functions, ranging from vector-based to structural kernels. KeLP allows to build complex kernel machine based systems, leveraging on the Java language and on a JSON interface to store and load classifiers configurations as well as to save the models to be reused.

For a deeper look, you can visit [What's inside KeLP page](#).

#### Downloads

KeLP is released under [Maven](#). To use it, please refer to the [Installation page](#). To download KeLP source code you can go to the github [KeLP page](#).

#### Authentication

[Log In](#)

#### News

- SAG's KeLP team ranked first at the [SemEval 2016 Community Question Answering Task](#) February 16, 2016
- [KeLP 2.0.2 released!](#) February 16, 2016
- [KeLP 2.0.1 released](#) January 13, 2016
- [The ECIR 2016 paper has been accepted!](#) December 30, 2015
- [KeLP 2.0.0 released](#) December 4, 2015
- [SAG with Reveal @ Maker Faire 2015, Rome!](#) October 16, 2015

<https://github.com/SAG-KeLP>

<http://sag.art.uniroma2.it/demo-software/kelp/>

# KELP applications: cQA

General Description | Subtasks | Data and Tools | Important Dates | **Results** | Call for Papers

## SemEval-2016 Task 3

**Task 3: Community Question Answering**

Building on the success of [SemEval 2015 Task 3](#) "Answer Selection in Community Question Answering" (see [the task description paper](#)), we propose an extension, which covers a full task on Community Question Answering (CQA) and which is, therefore, closer to a real application (see, e.g., [Qatar Living forum](#)).

CQA systems are gaining popularity online. Such systems are seldom moderated, quite open, and thus they have little restrictions, if any, on who can post and who can answer a question. On the positive side, this means that one can freely ask any question and expect some good, honest answers. On the negative side, it takes effort to go through all possible answers and to make sense of them. For example, it is not unusual for a question to have hundreds of answers, which makes it very time-consuming for the user to inspect and to winnow through them all. The present task could help to automate the process of finding good answers to new questions in a community-created discussion forum (e.g., by retrieving similar questions in the forum and by identifying the posts in the answer threads of those similar questions that answer the original question well).

In essence, the main CQA task can be defined as follows:

*"given (i) a new question and (ii) a large collection of question-comment threads created by a user community, rank the comments that are most useful for answering the new question"*

General Description | Subtasks | **Data and Tools**

## SemEval-2016 Task 3

**Results**

- ☐ The evaluation results can be found [here](#)
- ☐ The gold labels, submissions and scores for all teams can be found [here](#)
- ☐ The gold labels inside the test XML can be found [here](#)

**Task participants are strongly encouraged to submit a system description for SemEval 2016:**  
<http://alt.qcri.org/semeval2016/index.php?id=call-for-papers>

# KELP applications: cQA

[General Description](#)[Subtasks](#)[Data and Tools](#)[Important Dates](#)[Results](#)[Call for Papers](#)

## SemEval-2016 Task 3

### Results

- └ The evaluation results can be found [here](#)
- └ The gold labels, submissions and scores for all teams can be found [here](#)
- └ The gold labels inside the test XML can be found [here](#)

**Task participants are strongly encouraged to submit a system description paper by March 4, 2016:**

<http://alt.qcri.org/semEval2016/index.php?id=call-for-papers>

### Contact Info

#### Organizers

- ▶ Preslav Nakov, Qatar Computing Research Institute, HBKU
- ▶ Lluís Màrquez, Qatar Computing Research Institute, HBKU
- ▶ Alessandro Moschitti, Qatar Computing Research Institute, HBKU
- ▶ Walid Magdy, Qatar Computing Research Institute, HBKU
- ▶ James Glass, CSAIL-MIT
- ▶ Bilal Randeree, Qatar Living

**email :** *semEval-cqa@googlegroups.com*

### Other Info

#### Announcements



# KELP

## applications: cQA

General Description	Subtasks
<b>SemEval-2016 Task A</b>	

Team ID	Team Affiliation
ConvKN	Qatar Computing Research Institute,
ECNU	East China Normal University, China
ICL00	Institute of Computational Linguistics
ICRC-HIT	Intelligence Computing Research Center
ITNLP-AiKF	Intelligence Technology and Natural Language Processing
Kelp	University of Roma, Tor Vergata, Italy
MTE-NN	Qatar Computing Research Institute,
overfitting	University of Waterloo, Canada
PMI-cool	Sofia University, Bulgaria
QAIIIIT	IIIT Hyderabad, India
QU-IR	Qatar University, Qatar
RDI.team	RDI Egypt, Cairo University, Egypt
SemanticZ	Sofia University, Bulgaria
SLS	MIT Computer Science and Artificial Intelligence Laboratory
SUPer.team	Sofia University, Bulgaria; Qatar Computing Research Institute
UH-PRHLT	Pattern Recognition and Human Language Understanding
UniMelb	The University of Melbourne, Australia
UPC_USMBA	Universitat Politècnica de Catalunya

Table 5: The participating teams.

Submission	MAP	AvgRec	MRR	P	R	F1	Acc
<b>1 Kelp-primary</b>	<b>79.19<sub>1</sub></b>	<b>88.82<sub>1</sub></b>	<b>86.42<sub>1</sub></b>	<b>76.96<sub>1</sub></b>	<b>55.30<sub>8</sub></b>	<b>64.36<sub>5</sub></b>	<b>75.11<sub>2</sub></b>
ConvKN-contrastive1	78.71	88.98	86.15	77.78	53.72	63.55	74.95
SUPer.team-contrastive1	77.68	88.06	84.76	75.59	55.00	63.68	74.50
<b>2 ConvKN-primary</b>	<b>77.66<sub>2</sub></b>	<b>88.05<sub>3</sub></b>	<b>84.93<sub>4</sub></b>	<b>75.56<sub>2</sub></b>	<b>58.84<sub>6</sub></b>	<b>66.16<sub>2</sub></b>	<b>75.54<sub>1</sub></b>
<b>3 SemanticZ-primary</b>	<b>77.58<sub>3</sub></b>	<b>88.14<sub>2</sub></b>	<b>85.21<sub>2</sub></b>	<b>74.13<sub>4</sub></b>	<b>53.05<sub>10</sub></b>	<b>61.84<sub>8</sub></b>	<b>73.39<sub>5</sub></b>
ConvKN-contrastive2	77.29	87.77	85.03	74.74	59.67	66.36	75.41
<b>4 ECNU-primary</b>	<b>77.28<sub>4</sub></b>	<b>87.52<sub>5</sub></b>	<b>84.09<sub>6</sub></b>	<b>70.46<sub>6</sub></b>	<b>63.36<sub>4</sub></b>	<b>66.72<sub>1</sub></b>	<b>74.31<sub>4</sub></b>
SemanticZ-contrastive1	77.16	87.73	84.08	75.29	53.20	62.35	73.88
<b>5 SUPer.team-primary</b>	<b>77.16<sub>5</sub></b>	<b>87.98<sub>4</sub></b>	<b>84.69<sub>5</sub></b>	<b>74.43<sub>3</sub></b>	<b>56.73<sub>7</sub></b>	<b>64.39<sub>4</sub></b>	<b>74.50<sub>3</sub></b>
MTE-NN-contrastive2	76.98	86.98	85.50	58.71	70.28	63.97	67.83
SUPer.team-contrastive2	76.97	87.89	84.58	74.31	56.36	64.10	74.34
MTE-NN-contrastive1	76.86	87.03	84.36	55.84	77.35	64.86	65.93
SLS-contrastive2	76.71	87.17	84.38	59.45	67.95	63.41	68.13
SLS-contrastive1	76.46	87.47	83.27	60.09	69.68	64.53	68.87
<b>6 MTE-NN-primary</b>	<b>76.44<sub>6</sub></b>	<b>86.74<sub>7</sub></b>	<b>84.97<sub>3</sub></b>	<b>56.28<sub>9</sub></b>	<b>76.22<sub>1</sub></b>	<b>64.75<sub>3</sub></b>	<b>66.27<sub>8</sub></b>
<b>7 SLS-primary</b>	<b>76.33<sub>7</sub></b>	<b>87.30<sub>6</sub></b>	<b>82.99<sub>7</sub></b>	<b>60.36<sub>8</sub></b>	<b>67.72<sub>3</sub></b>	<b>63.83<sub>6</sub></b>	<b>68.81<sub>7</sub></b>
ECNU-contrastive2	75.71	86.14	82.53	63.60	66.67	65.10	70.95
SemanticZ-contrastive2	75.41	86.51	82.52	73.19	50.11	59.49	72.26
ICRC-HIT-contrastive1	73.34	84.81	79.65	63.43	69.30	66.24	71.28
<b>8 ITNLP-AiKF-primary</b>	<b>71.52<sub>8</sub></b>	<b>82.67<sub>9</sub></b>	<b>80.26<sub>8</sub></b>	<b>73.18<sub>5</sub></b>	<b>19.71<sub>12</sub></b>	<b>31.06<sub>12</sub></b>	<b>64.43<sub>9</sub></b>
ECNU-contrastive1	71.34	83.39	78.62	66.95	41.31	51.09	67.86
<b>9 ICRC-HIT-primary</b>	<b>70.90<sub>9</sub></b>	<b>83.36<sub>8</sub></b>	<b>77.38<sub>10</sub></b>	<b>62.48<sub>7</sub></b>	<b>62.53<sub>5</sub></b>	<b>62.50<sub>7</sub></b>	<b>69.51<sub>6</sub></b>
<b>10 PMI-cool-primary</b>	<b>68.79<sub>10</sub></b>	<b>79.94<sub>10</sub></b>	<b>80.00<sub>9</sub></b>	<b>47.81<sub>12</sub></b>	<b>70.58<sub>2</sub></b>	<b>57.00<sub>9</sub></b>	<b>56.73<sub>12</sub></b>
UH-PRHLT-contrastive1	67.57	79.50	77.08	54.10	50.11	52.03	62.45
<b>11 UH-PRHLT-primary</b>	<b>67.42<sub>11</sub></b>	<b>79.38<sub>11</sub></b>	<b>76.97<sub>11</sub></b>	<b>55.64<sub>10</sub></b>	<b>46.80<sub>11</sub></b>	<b>50.84<sub>11</sub></b>	<b>63.21<sub>10</sub></b>
UH-PRHLT-contrastive2	67.33	79.34	76.73	54.97	49.13	51.89	62.97
<b>12 QAIIIIT-primary</b>	<b>62.24<sub>12</sub></b>	<b>75.41<sub>12</sub></b>	<b>70.58<sub>12</sub></b>	<b>50.28<sub>11</sub></b>	<b>53.50<sub>9</sub></b>	<b>51.84<sub>10</sub></b>	<b>59.60<sub>11</sub></b>
QAIIIIT-contrastive2	61.93	75.22	69.95	49.48	49.96	49.72	58.93
QAIIIIT-contrastive1	61.80	75.12	69.76	49.85	50.94	50.39	59.24
Baseline 1 (IR)	<b>59.53</b>	<b>72.60</b>	<b>67.83</b>	—	—	—	—
Baseline 2 (random)	52.80	66.52	58.71	40.56	74.57	52.55	45.26
Baseline 3 (all 'true')	—	—	—	40.64	100.00	<b>57.80</b>	40.64
Baseline 4 (all 'false')	—	—	—	—	—	—	<b>59.36</b>

Table 1: **Subtask A, English (Question-Comment Similarity):** results for all submissions. The first column shows the rank of the primary runs with respect to the official MAP score. The second column contains the team’s name and its submission type (primary vs. contrastive). The following columns show the results for the primary, and then for other, unofficial evaluation measures. The subindices show the rank of the primary runs with respect to the evaluation measure in the respective column.

# ... a further example

- EvalIta 2018:

- IronIta: [Irony Detection in Italian Tweets \(IronITA\)](#)

- Task B: irony type



## 5.2 Task B: Different types of irony

- 34 The best performing UNITOR team is also the only team that participated to Task B with an unconstrained run.

Table 5: Results Task B. Unconstrained runs are marked by grey background

team name	id	F1-score			
		not-iro	iro	sarc	macro
UNITOR	2	0.668	0.447	0.446	0.520
UNITOR	1	0.662	0.432	0.459	0.518
ItaliaNLP	1	0.707	0.432	0.409	0.516
ItaliaNLP	2	0.693	0.423	0.392	0.503
Aspie96	1	0.668	0.438	0.289	0.465
<i>baseline-random</i>		0.503	0.266	0.242	0.337
venses-itgetarun	1	0.431	0.260	0.018	0.236
<i>baseline-mfc</i>		0.668	0.000	0.000	0.223
venses-itgetarun	2	0.413	0.183	0.000	0.199

- 35 All participating systems show an improvement over the baselines, with the exception of the only unsupervised system (venses-itgetarun, see details in Section 6).

# References

- Mitchell, Tom. M. 1997. *Machine Learning*. New York: McGraw-Hill.
- [Kernel machines, neural networks and graphical models](#), P. Frasconi, A. Sperduti, A. Starita, Rivista AI\*IA Numero speciale per i “50 anni di IA”, 2007.
- Very good video lectures by Andrew Ng (Stanford) <http://academicearth.org/courses/machine-learning>