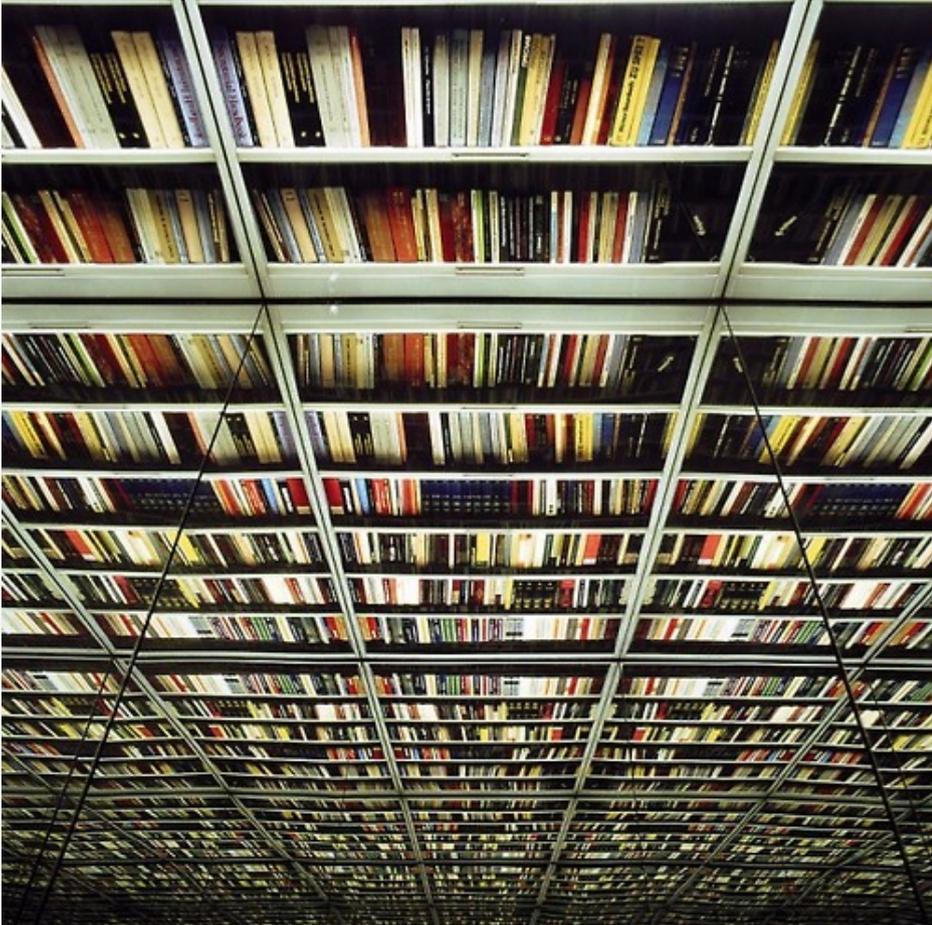


# Natural Language Processing and Information Retrieval - Major Tasks and Applications

Danas Zuokas  
Vilnius University

# The Problem 1



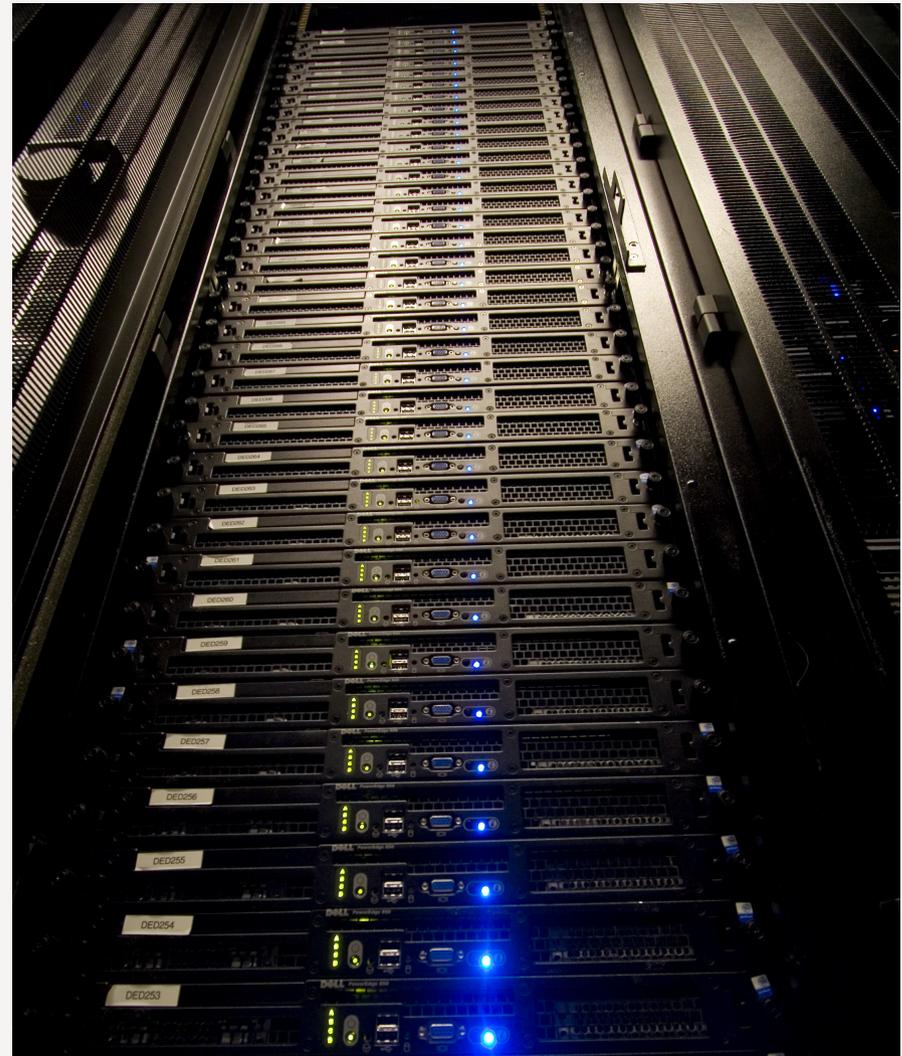
- The amount of data produced and stored keeps growing exponentially.
- Majority of that data is available in an semi- and unstructured form (text).
- Textual data contains significant business value.
- Effectively analyzing this type of data is of paramount importance.
- Traditional Information Technology techniques are not sufficient.

# The Solution

Machine Learning  
Text Analytics  
Information Retrieval Computational Linguistics  
Named Entity Recognition  
Natural Language Processing  
Sentiment Analysis

# The Problem 2

- The volume of data that is being employed by text analytics methods is vast.
- This requires a huge storage and computing infrastructure.
- It is very complex and expensive for a separate user to implement and maintain such systems.
- Actors with the advanced state-of-the-art scientific methodology knowledge and computing infrastructure are capable of providing text analytics service.



# The Solution

parallel computing

**grid computing**

computer cluster server farm

**cloud computing**

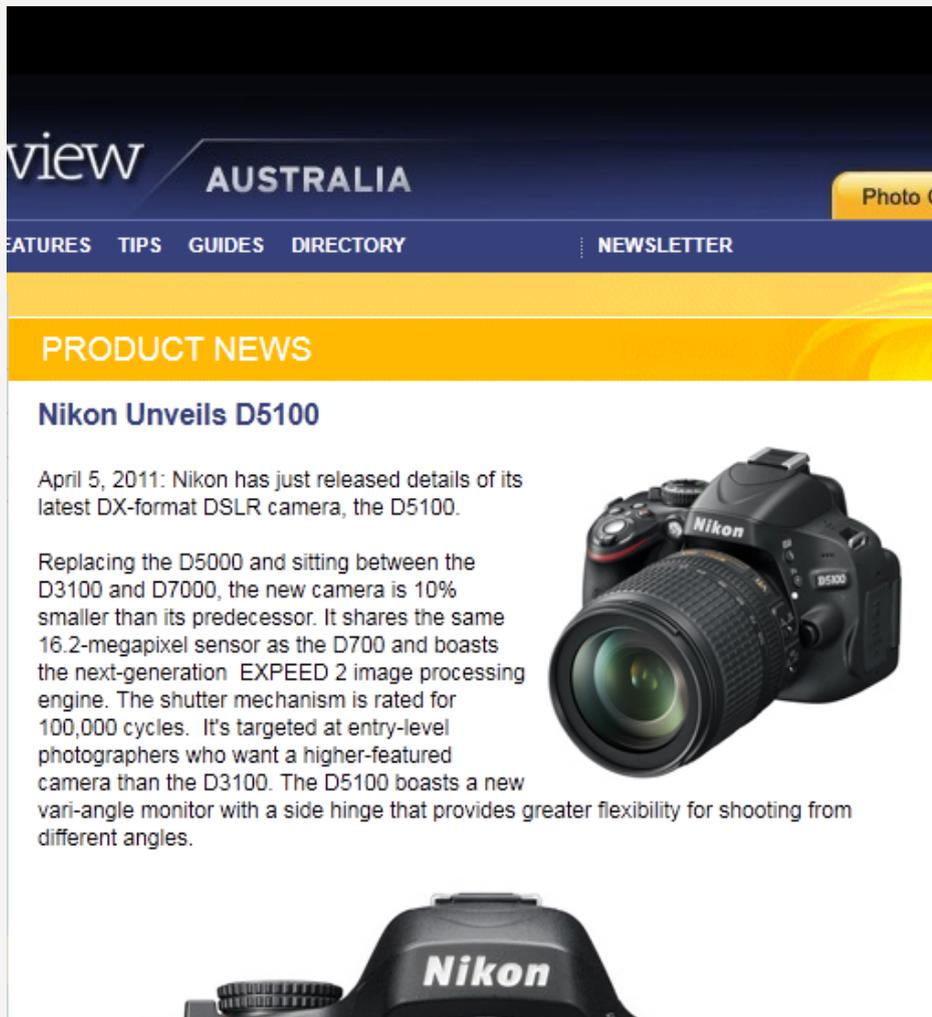
**supercomputer**

distributed computing

# A Model System

- Consider a system which deals with **events** as they are reported on the Web.
- The event is a segment of information describing **who does what and when**.
- On the one side a **service provider** would create a data pool of events.
- On the other side a user (**SME, public or other organization**) would get the **query-relevant events** with additional service:
  - **credibility** – a measure of importance of the event,
  - **sentiment** – opinions gathered and analyzed after the event was reported,
  - **associations** – a network of related events for further exploration.

# 2 Examples



view AUSTRALIA Photo G

FEATURES TIPS GUIDES DIRECTORY NEWSLETTER

## PRODUCT NEWS

### Nikon Unveils D5100

April 5, 2011: Nikon has just released details of its latest DX-format DSLR camera, the D5100.

Replacing the D5000 and sitting between the D3100 and D7000, the new camera is 10% smaller than its predecessor. It shares the same 16.2-megapixel sensor as the D700 and boasts the next-generation EXPEED 2 image processing engine. The shutter mechanism is rated for 100,000 cycles. It's targeted at entry-level photographers who want a higher-featured camera than the D3100. The D5100 boasts a new vari-angle monitor with a side hinge that provides greater flexibility for shooting from different angles.



washingtonpost.com > World

» THIS STORY: [READ +](#) | [WATCH +](#) | [Comments](#)

## Obama tells troops in Afghanistan: Success is within reach

By [Ernesto Londoño and Perry Bacon Jr.](#)  
Washington Post Foreign Service  
Friday, December 3, 2010; 8:54 PM

### GALLERY



#### President Obama makes surprise visit to Afghanistan

On an unannounced visit, Obama praised U.S. troops Friday and said his administration's surge over the past year has weakened the Taliban.

» [LAUNCH PHOTO GALLERY](#)

**THIS STORY**

- » Obama tells troops in Afghanistan: Success is within reach
- » [President Obama makes surprise visit to Afghanistan](#)
- » [Obama, in Afghanistan, says U.S. will succeed](#)

Obama landed almost exactly a year to the day after he authorized a 30,000-troop surge, and as his advisers finalize a comprehensive assessment of whether the expanded U.S. military presence is working.

With few decisive signs of progress to show at a time when many Americans are growing impatient with the war, the administration is once again reexamining its strategy to marginalize extremist groups and prop up the anemic Afghan government.

Network News [PROFILE](#)

[Recommend](#)

[Tweet](#) 0

[View More Activity](#)

**TOOLBOX**

[Resize](#) [Print](#) [E-mail](#)

# Named Entity Recognition



## DEMOS

### Named Entity Recognition Demo Results

The Named Entity Recognizer has identified the following named entities.

[**PER Obama**] tells troops in [**MISC Afghanistan: Success**] is within reach  
By [**PER Ernesto Londono**] and [**PER Perry Bacon Jr**] .  
[**ORG Washington Post Foreign Service**]  
Friday, December 3, 2010; 8:54 PM  
KABUL - President [**PER Obama**] told [**LOC U.S.**] troops in [**LOC Afghanistan**] that success is within reach during a brief, unexpected visit Friday night that came amid a series of embarrassing setbacks in the effort to turn around the faltering war.  
[**PER Obama**] landed almost exactly a year to the day after he authorized a 30,000-troop surge, and as his advisers finalize a comprehensive assessment of whether the expanded [**LOC U.S.**] military presence is working.  
With few decisive signs of progress to show at a time when many [**MISC Americans**] are growing impatient with the war, the administration is once again reexamining its strategy to marginalize extremist groups and prop up the anemic [**MISC Afghan**] government.

#### Key:

- **PER** - Person
- **ORG** - Organization
- **LOC** - Location
- **MISC** - Miscellaneous

# Event Extraction

**When:** April 5, 2011

**Who:** Nikon

**Did what:** has just released details

**About what:** of its latest DX-format DSLR camera, the D5100.

**Who:** President Obama

**Did what:** makes surprise visit

**Where:** to Afghanistan

**When:** at Friday, December 3, 2010

# Sentiment Analysis

04-05-2011 04:02 PM

#7

KmH ◦

TPF Junkie!



<b>Join Date:</b>	Apr 2009
<b>Location:</b>	Iowa
<b>Posts:</b>	12,800
<b>My Gallery:</b>	(0)
<b>My Photos Are OK to Edit</b>	
<b>Liked:</b>	179 times

The **D5100** has a lower quality AA filter than the D7000 has.

Not only does the D7000 have better high ISO performance, it will also have better image quality than the **D5100**.

Interestingly you failed to mention two of the D7000's most important features, it's Multi-CAM 4800 DX, 39 focus point (9 cross-type) auto focus module, and it's new 2016 pixel, 3D, RGB metering sensor.

The **D5100** is still using the 420 pixel RGB metering sensor, and the Multi-CAM 1000, 11 focus point (1 cross-type point) auto focus module.

Share

. .Keith. . . .

Check out my TPF - [Spring Cleaning Sale](#) - Nikon, Sekonic, Wacom, X-RITE (Only the Nikon stuff is still available, the rest

04-05-2011 07:48 PM

#10

blackxthink ◦

TPF Noob!



<b>Join Date:</b>	Apr 2011
<b>Location:</b>	Quebec, canada
<b>Posts:</b>	10
<b>My Gallery:</b>	(0)
<b>My Photos Are NOT OK to Edit</b>	
<b>Liked:</b>	0 times

d7000 + tokina 11-16mm 2.8 for extra wide

Depends on how much money you wanna spend.

The better gear you get at start , it will take longer for you to upgrade it.

I'd still get d7000 over **d5100** Only because it's a higher end camera.

5

But if you have a really restricted budget , id get **d5100** + cheaper lens

Share

# Event Grouping

Obama landed almost exactly a year to the day after he authorized a 30,000-troop surge, and as his advisers finalize a comprehensive assessment of whether the expanded U.S. military presence is working.

the Taliban.  
» LAUNCH PHOTO GAL

Friday's visit followed by a day the release of U.S. diplomatic cables that chronicled rampant corruption in the Afghan government and the toll that has taken during a war that is in its 10th year.

The U.S. effort has been marred by reversals and uncomfortable revelations in recent weeks. Last month, Afghan President Hamid Karzai angered U.S. officials after he said in an interview that U.S. troops were being overly aggressive in their pursuit of the Taliban.

The weather in and around Kabul was slightly hazy Friday afternoon, but it appeared to clear up by nightfall, when the president landed. Karzai and Obama met for roughly an hour last month in Lisbon during a NATO summit.

The White House expects to complete the process by Dec. 18, before Obama travels to Hawaii for the holidays.

# Faceted Search

**NY TIMES ARTICLE SEARCH**

DESCRIPTION	YEAR	LOCATION	PERSON	ORGANIZATION
afghanistan war (2001- ) 84	2009 70	afghanistan 130	karzai, hamid 130	taliban 56
united states international 67	2008 7	pakistan 12	obama, barack 56	north atlantic treaty organ 19
united states defense and 54		united states 8	abdullah, abdullah 18	al qaeda 8
elections 39		kandahar (afghanistan) 5	petraeus, david h 8	state department 7
terrorism 30		marja (afghanistan) 4	mcchrysal, stanley a 7	wikileaks 5
politics and government 30		kabul (afghanistan) 4	clinton, hillary rodham 6	united nations 5
ethics 19	2010 50	iraq 2	bush, george w 6	central intelligence agency 4
editorials 13	2011 3	iran 2	karzai, ahmed wali 5	defense department 2
frauds and swindling 12		philippines 1	zardari, asif ali 4	international monetary fund 1
defense and military force 10		munich (germany) 1	holbrooke, richard c 4	international committee of 1
united states armament an 8		helmand province (afghani 1	ghani, ashraf 3	house of representatives 1
voting and voters 7		great britain 1	gates, robert m 3	haqqani network 1
drug abuse and traffic 5		germany 1	fahim, muhammad 3	der spiegel 1
		farah province (afghanista 1	biden, joseph r jr 3	democratic party 1
		europa 1	atmar, muhammad hanif 2	asian development bank 1

## 130 RESULTS

<p><b>NATO Helicopters Kill 9 Afghan Boys Collecting Firewood for</b></p> <p>2011/03/03</p>	<p><b>THE SHADOW WAR; A FORMER SPY, NOW OPERATING HIS OWN C.I.A.</b></p> <p>2011/01/23</p>	<p><b>EDITORIAL; President Karzai's Latest</b></p> <p>2011/01/21</p> <p>Editorial scores Afghanistan Pres Hamid Karzai for ordering Parliament to delay its opening session while an unconstitutional court he appointed re-investigates charges of fraud in the 2010 parliamentary elections; urges Pres Obama to make it clear to Karzai publicly and privately that he is accountable to his people and his</p>
---	--	--

# Impact

The user would benefit in saving operational costs when exploring web content. The service would benefit SME's, public and other organizations. To give a few examples:

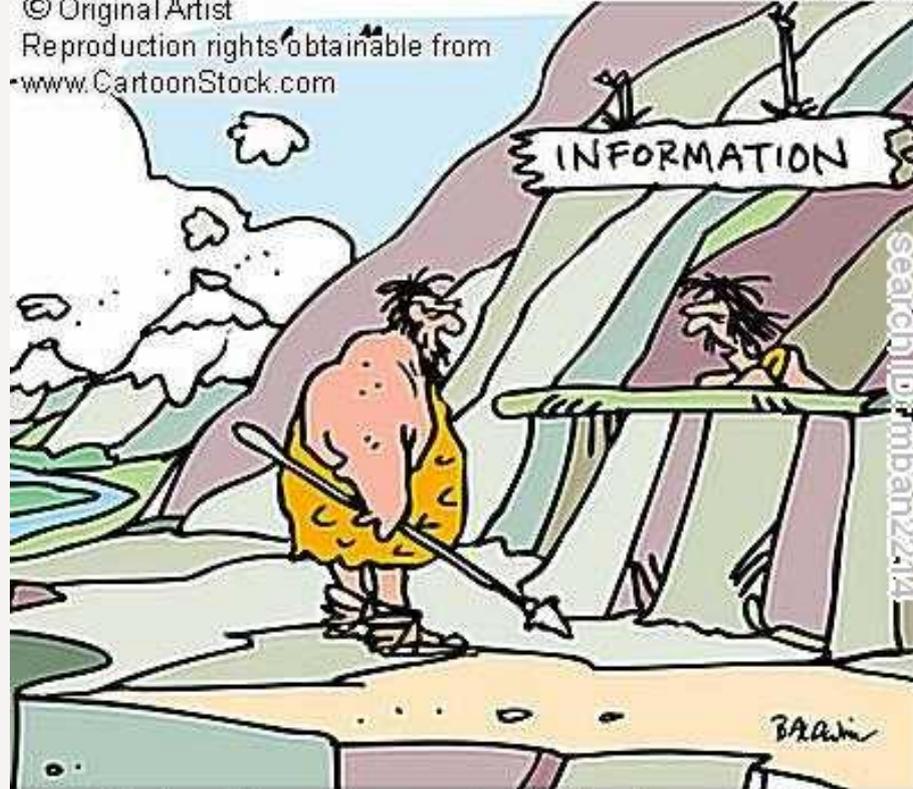
- **Content and news providers** would use data pool to add functionality to their publications.
- The tool would help **governmental bodies** to evaluate the feedback from the citizens after the decisions are made.
- **Marketing/Public Relations companies** would use event pool to evaluate outcomes of their campaigns.
- **Predictive analytics companies** would use event pool to create predictive models (e.g. portfolio manager). Some events change sentiment on the market which then moves asset price.

# Conclusions

- **Similar services** that would use text analytics tools also are and will be developed.
- Such systems are very **complex and expensive** in order to be implemented and maintained by separate users on their own.
- Only actors with the knowledge of advanced state-of-the-art **scientific methodology and computing infrastructure** are capable of creating such systems.
- Users would identify the demand for **specific functionalities** in such a way determining “the shape” of a service.

© Original Artist  
Reproduction rights obtainable from  
www.CartoonStock.com

© Mike Baldwin / Corbis



"I don't have any yet. We just opened."