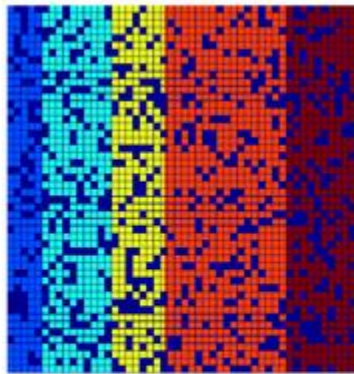# Recap: how to build such a space

- Solution
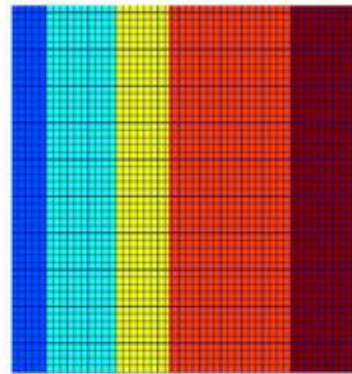  - Low rank matrix approximation

*Imagine this is \*true\* concept-document matrix*
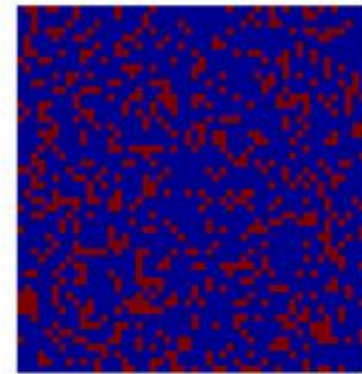


Matrix of corrupted observations → Underlying low-rank matrix + Sparse error matrix

*Imagine this is our observed term-document matrix*

*Random noise over the word selection in each document*

# Recap: Latent Semantic Analysis (LSA)

- Solve LSA by SVD

  *Map to a lower dimensional space*

  $$\hat{Z} = \underset{Z|rank(Z)=k}{\operatorname{argmin}} \|C - Z\|_F$$

  $$= \underset{Z|rank(Z)=k}{\operatorname{argmin}} \sqrt{\sum_{i=1}^{M} \sum_{j=1}^{N} \left(C_{ij} - Z_{ij}\right)^2}$$

  $$= C_{M \times N}^{k}$$

  – Procedure of LSA

  1. Perform SVD on document-term adjacency matrix
  2. Construct $C_{M \times N}^{k}$ by only keeping the largest $k$ singular values in $\Sigma$ non-zero

# Introduction to Natural Language Processing
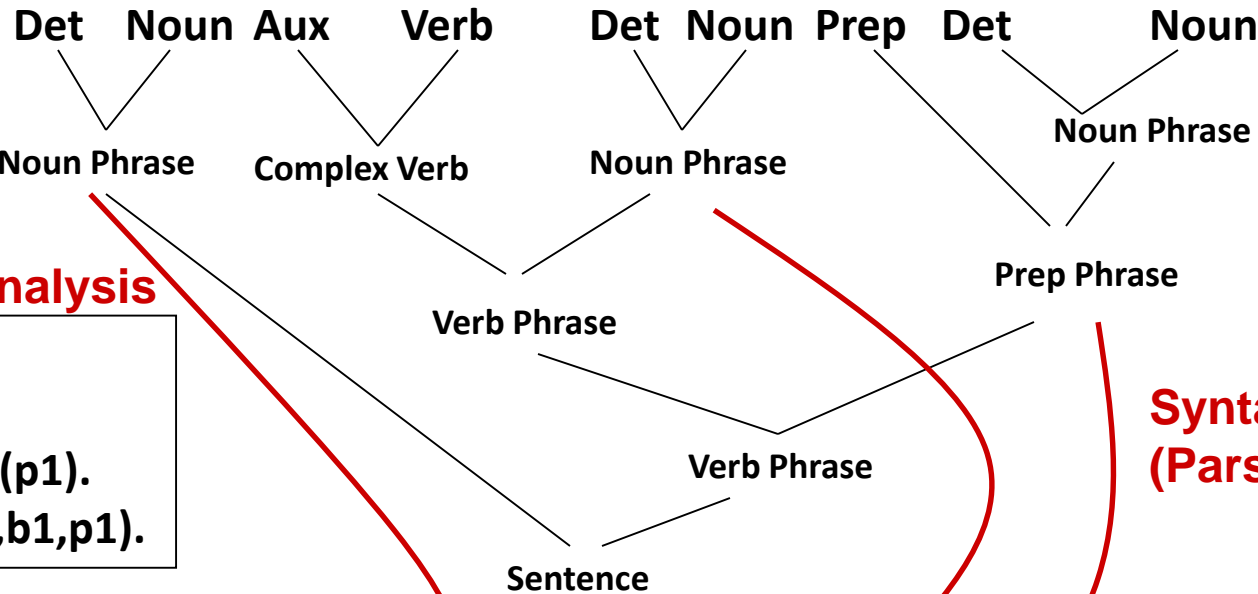
Hongning Wang

CS@UVa

# What is NLP?

**Arabic text**     كلب هو مطاردة صبي في الملعب.

How can a computer make **sense** out of this **string?**

| | |
|---|---|
| **Morphology** | - What are the basic units of meaning (words)?<br>- What is the meaning of each word? |
| **Syntax** | - How are words related with each other? |
| **Semantics** | - What is the "combined meaning" of words? |
| **Pragmatics** | - What is the "meta-meaning"? (speech act) |
| **Discourse** | - Handling a large chunk of text |
| **Inference** | - Making sense of everything |

# An example of NLP

A   dog   is   chasing   a   boy   on   the   playground.

**Det   Noun   Aux   Verb   Det   Noun   Prep   Det   Noun**

Noun Phrase   Complex Verb   Noun Phrase   Noun Phrase

**Lexical analysis (part-of-speech tagging)**

Verb Phrase

Prep Phrase

**Semantic analysis**

Dog(d1).
Boy(b1).
Playground(p1).
Chasing(d1,b1,p1).

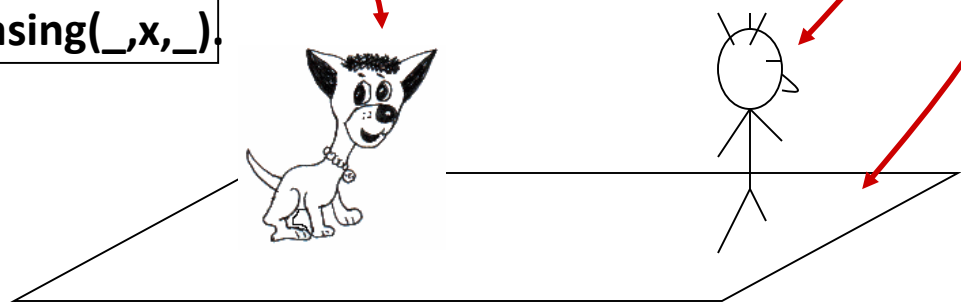Verb Phrase

**Syntactic analysis (Parsing)**

+

Scared(x) if Chasing(_,x,_).

Sentence

Scared(b1)

**Inference**

**A person saying this may be reminding another person to get the dog back...**

**Pragmatic analysis (speech act)**

CS@UVa                    CS6501: Text Mining                    5

- *Automatically answer our emails*
- *Translate languages accurately*
- *Help us manage, summarize, and aggregate information*
- *Use speech as a UI (when needed)*
- *Talk to us / listen to us*

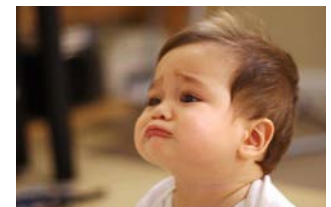# If we can do this for all the sentences in all languages, then …
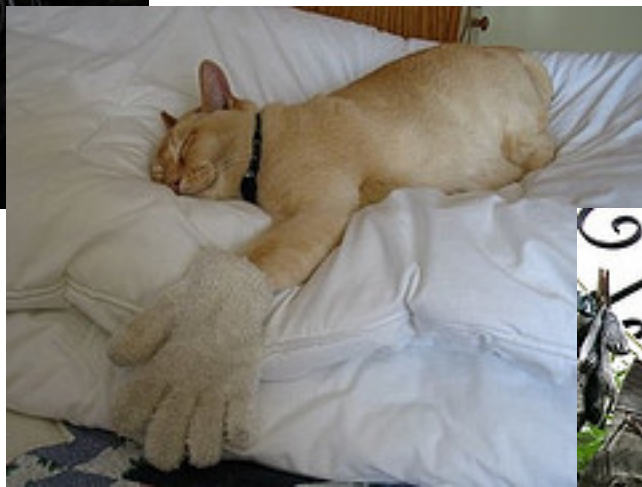
**BAD N**

- **Unf** **ght now.**
- **Gen** **"**

# NLP is difficult!!!!!!!

- Natural language is designed to make human communication efficient. Therefore,
  - We omit a lot of "common sense" knowledge, which we assume the hearer/reader possesses
  - We keep a lot of ambiguities, which we assume the hearer/reader knows how to resolve

- This makes EVERY step in NLP hard
  - Ambiguity is a "killer"!
  - Common sense reasoning is pre-required

CS6501: Text Mining

# An example of ambiguity

- Get the cat with the gloves.

# Examples of challenges

- Word-level ambiguity
  - "design" can be a noun or a verb (Ambiguous POS)
  - "root" has multiple meanings (Ambiguous sense)
- Syntactic ambiguity
  - "natural language processing" (Modification)
  - "A man saw a boy _with a telescope_." (PP Attachment)
- Anaphora resolution
  - "John persuaded Bill to buy a TV for himself." (himself = John or Bill?)
- Presupposition
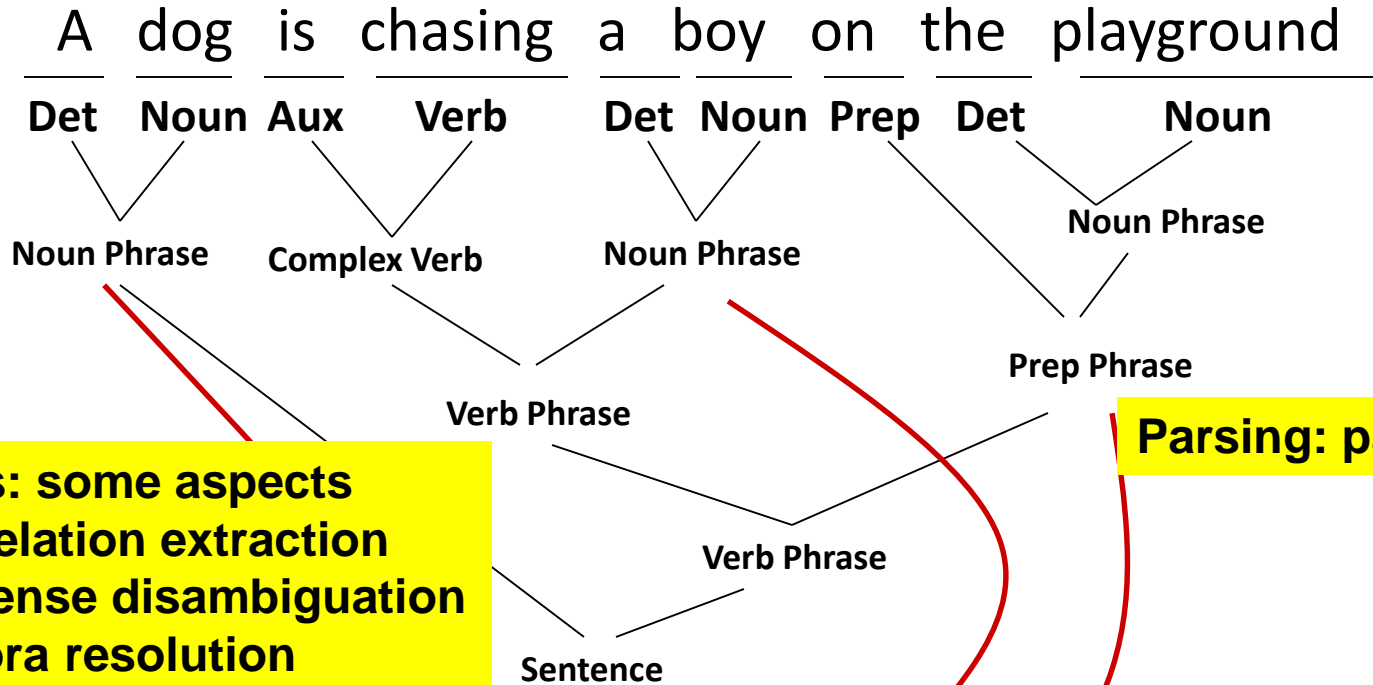  - "He has quit smoking." implies that he smoked before.

Despite all the challenges, research in NLP has also made a lot of progress…

# A brief history of NLP

- Early enthusiasm (1950's): Machine Translation
  - Too ambitious
  - Bar-Hillel report (1960) concluded that fully-automatic high-quality translation could not be accomplished without knowledge  (Dictionary + Encyclopedia)
- Less ambitious applications (late 1960's & early 1970's): Limited success, failed to scale up
  - Speech recognition
  - Dialogue (Eliza)   **Shallow understanding**          **Deep understanding in limited domain**
  - Inference and domain knowledge (SHRDLU="block world")
- Real world evaluation (late 1970's – now)
  - Story understanding (late 1970's & early 1980's)   **Knowledge representation**
  - Large scale evaluation of speech recognition, text retrieval, information extraction (1980 – now)   **Robust component techniques**
  - Statistical approaches enjoy more success (first in speech recognition & retrieval, later others)   **Statistical language models**
- Current trend:
  - Boundary between statistical and symbolic approaches is disappearing.
  - We need to use all the available knowledge   **Applications**
  - Application-driven NLP research (bioinformatics, Web, Question answering...)
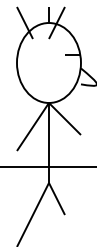
# The state of the art

A   dog   is   chasing   a   boy   on   the   playground

**Det**   **Noun**   **Aux**   **Verb**   **Det**   **Noun**   **Prep**   **Det**   **Noun**

**POS Tagging: 97%**

**Noun Phrase**   **Complex Verb**   **Noun Phrase**   **Noun Phrase**

**Prep Phrase**

**Verb Phrase**

**Semantics: some aspects**
- **Entity/relation extraction**
- **Word sense disambiguation**
- **Anaphora resolution**

**Parsing: partial >90%**

**Verb Phrase**

**Sentence**

**Inference: ???**

**Speech act analysis: ???**

# Machine translation

# Machine translation

# Dialog systems



Apple's siri system



Google search

# Information extraction



Interstellar (2014)

PG-13 · 2hr 49min · Science Fiction

IMDb                                    8.9/10 ★★★★☆
Rotten Tomatoes                         73% ★★★☆☆

In the near future around the American Midwest, Cooper an ex-science engineer and pilot, is tied to his farming land with his daughter Murph and son Tom. As devastating sandstorms ravage earths crops, the people of Earth realize their life here … +

en.wikipedia.org

Boxoffice gross: $779 million USD
Estimated budget: $165 million USD
Release date: Nov 05, 2014
Director: Christopher Nolan
Screenwriters: Christopher Nolan · Jonathan Nolan
Music by: Hans Zimmer

Watch movie
▶ Watch trailer on YouTube

Cast                                    See all (20+)

Matthew McConaug… Cooper | Anne Hathaway Brand | Jessica Chastain Murph | Casey Affleck | Wes Bentley Doyle

University of Virginia

| Established | 1819 |
| Type | Public Flagship |
| Endowment | US$6.4 billion[1] |
| Budget | US$2.7 billion (2013—excludes capital spending) |
| President | Teresa A. Sullivan |
| Academic staff | 2,102 |
| Undergraduates | 14,898[2] |
| Postgraduates | 6,340[2] |
| Location | Charlottesville, Virginia, United States |
| Campus | Suburban 1,682 acres (6.81 km$^2$) |

Google Knowledge Graph

Wiki Info Box

# Information extraction

Search: [_____] [eng ▼]                    **<Albert_Einstein>**

← <Elsa_Einstein>                                                    "albert. ainctain"@jbo
← <Mileva_Marić>      <isMarriedTo>                                  "Albert Einstein"@afr

## Recently-Learned Facts  [twitter]                                    [Refresh]

| instance | iteration | date learned | confidence | | |
|----------|-----------|--------------|------------|---|---|
| tear_drop_tomatoes is an agricultural product | 887 | 27-nov-2014 | 93.1 | 👍 | 👎 |
| ryan_mckenzie is a professor | 886 | 21-nov-2014 | 90.2 | 👍 | 👎 |
| fiorina_161 is a planet | 889 | 07-dec-2014 | 92.8 | 👍 | 👎 |
| critical_thinking_in_health_science is a cognitive action | 886 | 21-nov-2014 | 99.0 | 👍 | 👎 |
| fateful_new_year is a monarch | 886 | 21-nov-2014 | 99.0 | 👍 | 👎 |
| tony_martin has been charged with murder | 890 | 11-dec-2014 | 100.0 | 👍 | 👎 |
| sen__joe_biden is a U.S. politician who holds the office of vice_president | 887 | 27-nov-2014 | 93.8 | 👍 | 👎 |
| hat is a clothing item to go with blue_jeans | 889 | 07-dec-2014 | 93.8 | 👍 | 👎 |
| statistics is headquartered in the country the_usa | 891 | 18-dec-2014 | 98.4 | 👍 | 👎 |
| eoin_colfer wrote the book artemis_fowl | 886 | 21-nov-2014 | 100.0 | 👍 | 👎 |

<Abelardo_[

<Abraham_Pais>
<Abram_L._Sachar>          **CMU Never-Ending Language Learning**
<Absent-minded_professor>
<Absorption_refrigerator>

<Albert_Einstein's_brain>
<Alfred_Kleiner>
...alen_der_Physik>
<Annus_Mirabilis_papers>

896–1954)>

## YAGO Knowledge Base

# Building a computer that 'understands' text: The NLP pipeline

# Tokenization/Segmentation

- Split text into words and sentences
  - Task: what is the most <span style="color:red">likely</span> segmentation /tokenization?

There was an earthquake near D.C. I've even felt it in Philadelphia, New York, etc.

There + was + an + earthquake + near + D.C.

I + ve + even + felt + it + in + Philadelphia, + New + York, + etc.

# Part-of-Speech tagging

- Marking up a word in a text (corpus) as corresponding to a particular part of speech
  - Task: what is the most likely tag sequence

A + dog + is + chasing + a + boy + on + the + playground

| A | dog | is | chasing | a | boy | on | the | playground |
|-----|------|-----|---------|-----|------|------|-----|------------|
| Det | Noun | Aux | Verb | Det | Noun | Prep | Det | Noun |

# Named entity recognition

- Determine text mapping to proper names
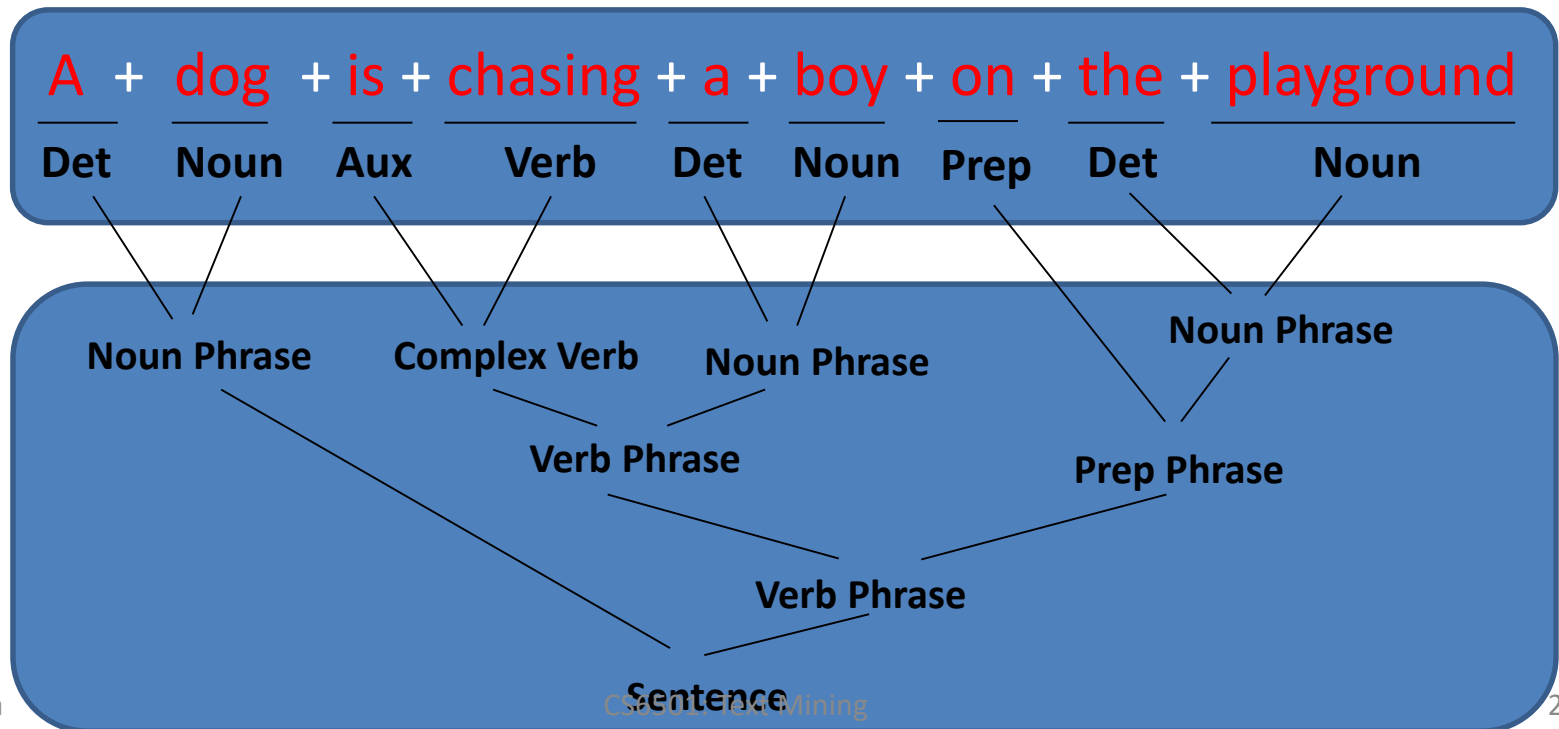  - Task: what is the most likely mapping

> Its initial Board of Visitors included U.S. Presidents Thomas Jefferson, James Madison, and James Monroe.

> Its initial Board of Visitors included U.S. Presidents Thomas Jefferson, James Madison, and James Monroe.

**Organization**, **Location**, **Person**

# Syntactic parsing

- Grammatical analysis of a given sentence, conforming to the rules of a formal grammar
  - Task: what is the most likely grammatical structure

A + dog + is + chasing + a + boy + on + the + playground

Det   Noun   Aux   Verb   Det   Noun   Prep   Det   Noun

Noun Phrase   Complex Verb   Noun Phrase

Noun Phrase

Verb Phrase

Prep Phrase

Verb Phrase

**Sentence**

# Relation extraction

- Identify the relationships among named entities
  - Shallow semantic analysis

Its initial Board of Visitors included U.S. Presidents Thomas Jefferson, James Madison, and James Monroe.

1. Thomas Jefferson Is_Member_Of Board of Visitors
2. Thomas Jefferson Is_President_Of U.S.

# Logic inference

- Convert chunks of text into more formal representations
    - Deep semantic analysis: e.g., first-order logic structures

Its initial Board of Visitors included U.S. Presidents Thomas Jefferson, James Madison, and James Monroe.

$\exists x$ (Is_Person($x$) & Is_President_Of($x$,'U.S.') & Is_Member_Of($x$,'Board of Visitors'))
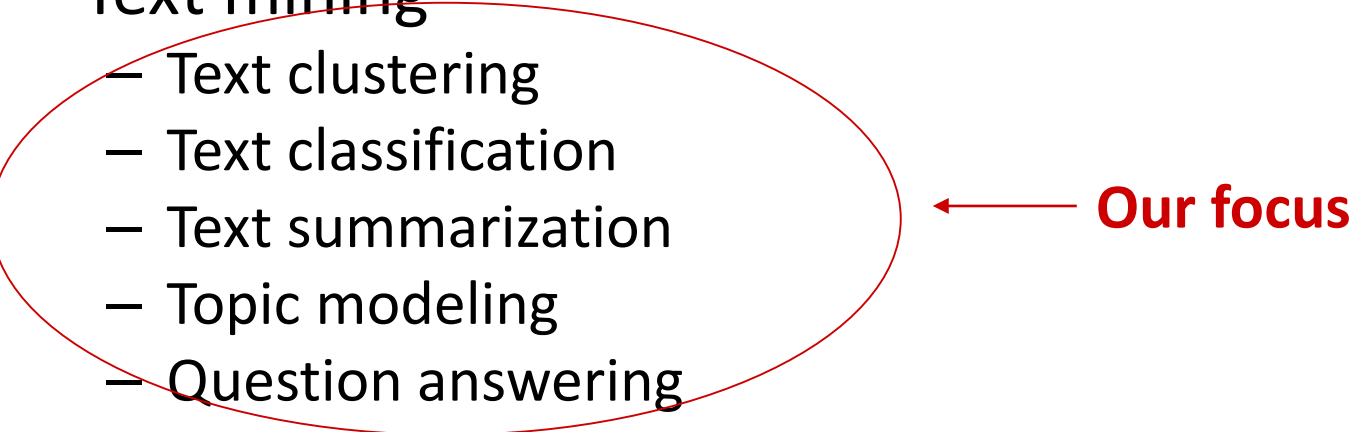
# Towards understanding of text

More than a decade ago, Carl Lewis stood on the threshold of what was to become the greatest athletics career in history. He had just broken two of the legendary Jesse Owens' college records, but never believed he would become a corporate icon, the focus of hundreds of millions of dollars in advertising. His sport was still nominally amateur. Eighteen Olympic and World Championship gold medals and 21 world records later, Lewis has become the richest man in the history of track and field -- a multi-millionaire.

- Who is Carl Lewis?

- Did Carl Lewis break any records?

# Major NLP applications

- Speech recognition: e.g., auto telephone call routing
- Text mining
  - Text clustering
  - Text classification
  - Text summarization          ← **Our focus**
  - Topic modeling
  - Question answering
- Language tutoring
  - Spelling/grammar correction
- Machine translation
  - Cross-language retrieval
  - Restricted natural language
- Natural language user interface

# NLP & text mining

- Better NLP => Better text mining

- Bad NLP => Bad text mining?

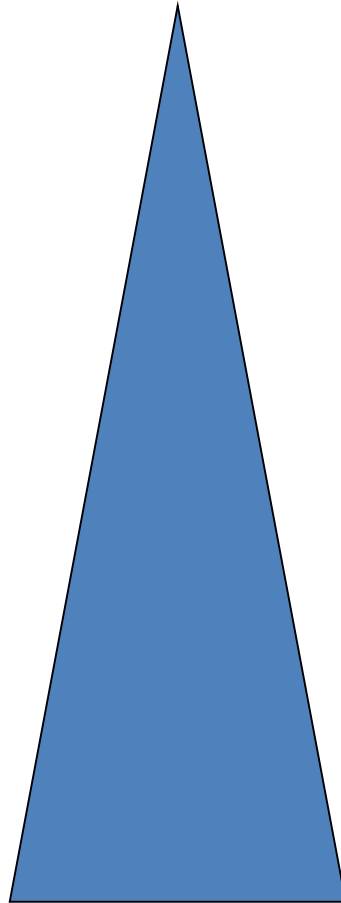**Robust, shallow NLP tends to be more useful than deep, but fragile NLP.**

**Errors in NLP can hurt text mining performance…**
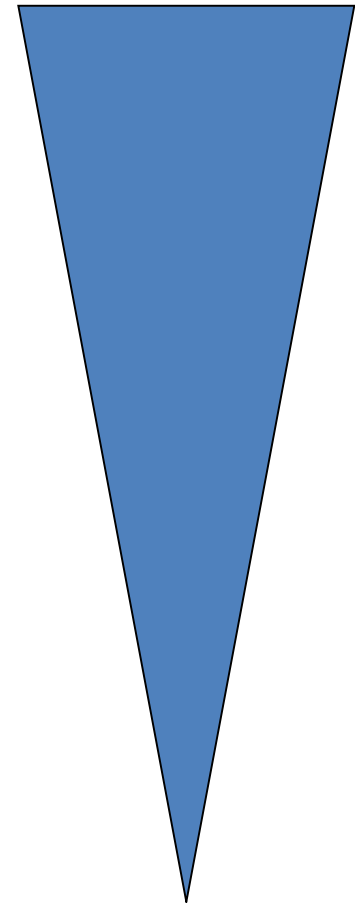
# How much NLP is really needed?

**Tasks**

**Dependency on NLP**

**Scalability**

Classification

Clustering

Summarization

Extraction

Topic modeling

Translation
Dialogue

Question
Answering

Inference
Speech Act

# So, what NLP techniques are the most useful for text mining?

- <u>Statistical NLP</u> in general.
- The need for high robustness and efficiency implies the dominant use of <u>simple models</u>

# What you should know

- Different levels of NLP

- Challenges in NLP

- NLP pipeline