

Natural Language Processing (COM4513/6513)

Week 5: Information Extraction

Prof. Jochen L. Leidner

[⟨leidner@acm.org⟩](mailto:leidner@acm.org)



The
University
Of
Sheffield.

University of Sheffield, Department of Computer Science
8 March 2018

Copyright ©2018 by Jochen L. Leidner. All rights reserved.

In this session, we aim to:

- get introduced to information extraction, its concepts and history
- learn about different approaches to IE

“Information Extraction” Defined

Information Extraction (IE) is

- the extraction of structured (relational) data from unstructured (= textual) sources
- a practically-motivated engineering discipline (models not necessarily inspired by nature)
- the use of natural language processing techniques to populate the slots of structured templates with appropriate fillers

IE was conceived as a shortcut to build useful systems when full text understanding based on syntactic parsing was beyond the state of the art.

“Concepcion, 23 Aug 88 (Santiago Domestic Service) – Police sources have reported that unidentified individuals planted a bomb in front of a Mormon Church in Talcahuano District. The bomb, which exploded and caused property damage worth 50,000 pesos, was placed at a chapel of the Church of Jesus Christ of Latter-Day Saints located at No 3856 Gomez Carreno Street.

The shock wave destroyed a wall, the roof, and the windows of the church, but did not cause any injuries.

Carabineros bomb squad personnel immediately went to the location and discovered that the bomb was made of 50 grams of an-fo ammonium nitrate-fuel oil blasting agents and a slow fuse.”

IE Example: Terrorist Attack Event

Concepcion, **23 Aug 88** “(Santiago Domestic Service) – Police sources have reported that **unidentified individuals** planted a **bomb** in front of a **Mormon Church** in **Talcahuano District**. The bomb, which exploded and caused **property damage worth 50,000 pesos**, was placed at a chapel of the Church of Jesus Christ of Latter-Day Saints located at No 3856 Gomez Carreno Street.

The shock wave destroyed a wall, the roof, and the windows of the church, but **did not cause any injuries**.

Carabineros bomb squad personnel immediately went to the location and discovered that the bomb was made of 50 grams of an-fo ammonium nitrate-fuel oil blasting agents and a slow fuse.”

— MUC4 story TST4-MUC4-0001

The “Terrorist Attack” Template

Terrorist Attack	
Type of Attack:	_____
Perpretator:	_____
Target:	_____
Location:	_____
Time:	_____
Casualties:	_____
Injured:	_____
Material Damage:	_____

The “Terrorist Attack” Template

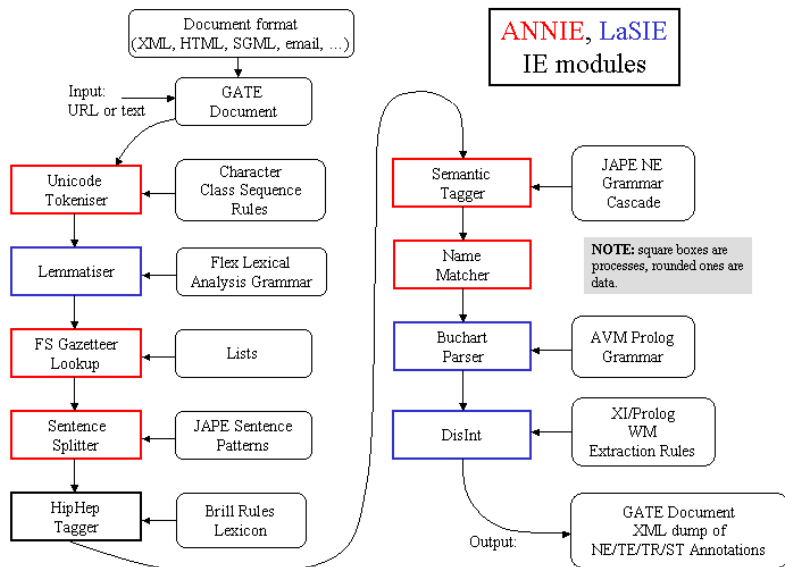
Terrorist Attack	
Type of Attack:	bomb (ATTACK>BOMBING)
Perpretator:	unidentified individuals (unknown)
Target:	a Mormon church
Location:	Talcahuano District (Chile>Talcahuano)
Time:	23 Aug 88 (1988-08-23 00:00:00)
Casualties:	_____ (0)
Injured:	did not cause any injuries (0)
Material Damage:	property worth 50,000 pesos (50000)

A Short History of Information Extraction

- 1981: NYU “Linguistic String” project (N. Sager)
- 1982: DeJong’s FRUMP system: ‘sketchy scripts’
- Carnegie Group build JASPER IE system for Reuters (Andersen et al., 1986)
- 1987/1989: MUCK I+II: Naval operations messages
- 1991-1998 MUC 3-7: Message Understanding Contest
- 2000-2004: ACE: Automatic Content Extraction
 - from text spans to abstract entities
 - English, Chinese, Arabic
- 2010s: First neural approaches to IE

- **Preprocessor/Tokenizer**: split story into units and ultimately word tokens
- **Gazetteer**: lexical look-up of important (to your task) words/phrases
- **POS tagger**: tag/disambiguate words w.r.t. parts of speech
- **Chunk parser**: find basic noun and verb phrases
- **Named entity tagger**: identify and classify proper names
- **Relationship tagger**: find relations between entities
- **Event Template Analyzer**: populate fact/event templates
- The result is a **structured fact/event database**

IE System Architecture – Example



- **Rule-based:** human experts (computational linguists) manually write general linguistic rules and task-specific extraction rules.
Trigger keywords, regular expressions, pattern/action, (cascaded) Finite State Transducers (FSTs)
- Supervised **Machine learning-based:** humans (domain experts) manually annotate text spans indicating entities, relations, facts etc. in a training corpus, and an expert (computational linguist) formulates a set of features; these get used to extract information if statistically correlated with classes of entities, relations etc. sought.
- Insight: shallow processing works well: SRI Tacitus → SRI FASTUS (Hobbs et al. 1992)

Example Patterns (FASTUS, Appelt et al., 1993)

killing of <HumanTarget>
<GovtOfficial> accused <PerpOrg>
bomb was placed by <Perp>
on <PhysicalTarget>
<Perp> attacked <HumanTarget>
<PhysicalTarget> with <Device>
<HumanTarget> was injured
<HumanTarget>'s body

From Entity over Relation to Scenario Template

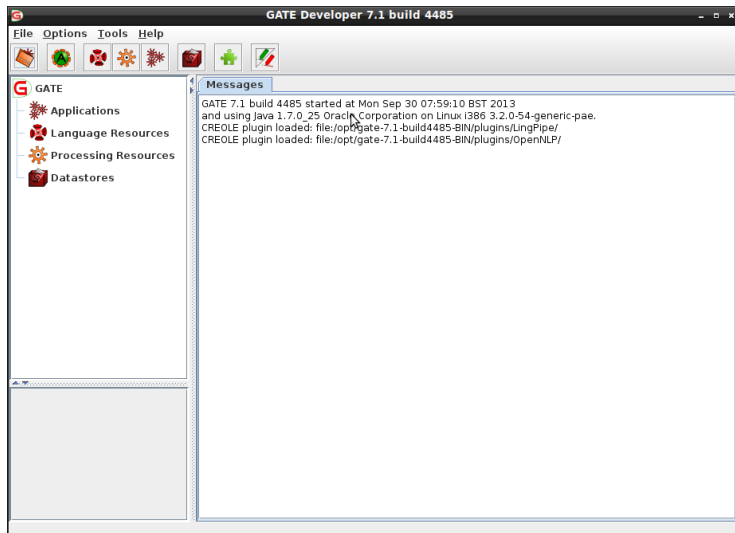
- **Template Element Recognition (TE)**: extract information pertaining to organizations, persons and artifacts (NE tagging, parsing)
- **Relationship Extraction (RE)**: extract information about how individual entities stand in relationship to each other, drawing on a pre-defined inventory of relation types
- **Scenario Template Recognition (ST)**: extract pre-specified event information and relate the event to particular organizations, persons and/or artifacts (slot fillers)
- Information for filling a template often often spread across several sentences

- Create a **gold data** set reference corpus with **ground truth**)
- Split gold data into three parts:
 - **development/training set**: used to study the data, train machine learning processes; can be inspected
 - **development test** (“devtest”) set: cannot be inspected, cannot be used for training; repeatedly used to measure improvements of system quality by comparing system output with ground truth.
 - **test set**: cannot be inspected; only used once for final evaluation run at project end. Completely **unseen data** (to the system and developers).
- Gold data split:
could be e.g. 80% train : 10% dev-test : 10% test

The GATE framework: From GUI to Guts

- GATE: General Architecture for Text Engineering - an open-source framework for *language engineering*
- Under development at the University of Sheffield for a long time (Cunningham et al., 2013)
- Current version: 8.4.2 (as of today) available from <http://gate.ac.uk>
- Java-based platform comprising GUI, APIs, and a workflow system
- GATE comes with pre-existing, contributed data (*Language Resources*) and components (*Processing Resources*)

Constructing IE Pipelines with GATE (1/2)



Constructing IE Pipelines with GATE (2/2)

The screenshot displays the GATE Developer 7.1 build 4485 interface. The left sidebar shows a tree view of resources, including Applications, Language Resources, Processing Resources, and a Corpus Pipeline_0003D. The main workspace is divided into several panels. The top panel shows the 'Corpus Pipeline...' configuration, with a 'Messages' tab selected. Below this, the 'Loaded Processing resources' panel lists 'IE' and 'RegEx Sentence Splitter_00017R'. The 'Selected Processing resources' panel lists 'Document Reset PR_00016', 'ANNIE Sentence Splitter_0003A', 'ANNIE English Tokeniser_00018', 'ANNIE Gazetteer_00042', and 'NE ANNIE NE Transducer_0001C'. The 'Corpus:' dropdown is set to 'GATE Corpus_00022'. The 'No processing resource selected...' message is displayed. The bottom panel shows the 'Serial Application Editor' and 'Initialisation Parameters' tabs. The 'Run this Application' button is visible.

GATE Developer 7.1 build 4485

File Options Tools Help

RegEx Sentence ... Document Reset ... Corpus Pipeline...

GATE Corpus_000... ANNIE NE Transd... ANNIE English T...

Messages IE test.txt_00021

Loaded Processing resources

Name
IE
RegEx Sentence Splitter_00017R

Selected Processing resources

Name
Document Reset PR_00016
ANNIE Sentence Splitter_0003A
ANNIE English Tokeniser_00018
ANNIE Gazetteer_00042
NE ANNIE NE Transducer_0001C

Corpus: GATE Corpus_00022

No processing resource selected...

Name	Type	Required	Value
------	------	----------	-------

Run this Application

Serial Application Editor Initialisation Parameters

- Part of the GATE platform
- Language to specify FSTs following the “patterns & actions” paradigm
- Sets of rules, broken down into processing phases
- Each rule tests matching conditions and typically adds annotations in the affirmative case

JAPE Rules – Regular Expressions over Annotations

Jape rules are organized into phases (cascades).

Each JAPE rule has three parts:

- **Header:** name of the rule
- **Left-Hand Side (LHS):** pattern
- **Right-Hand Side (RHS):** action – e.g. add annotations
- Structure:

```
Rule: TagUnknownName
(
  {Token.category == NNP}
):x
-->
  :x.Unknown = { kind = "PN", rule = TagUnknownName }
```

- See also: <http://gate.ac.uk/sale/tao/splitch8.html#chap:jape>

JAPE Rule Format – Example 1: URL Prefix

Phase: UrlPre

Input: Token SpaceToken

/* important: specify input! */

Options: control = appelt

Rule: Urlpre

```
( ({Token.string == "http"} |  
  {Token.string == "ftp"})  
  {Token.string == ":"}  
  {Token.string == "/"}  
  {Token.string == "/"}  
  ) |  
  ({Token.string == "www"}  
   {Token.string == "."})  
  )
```

):urlpre

-->

:urlpre.UrlPre = {rule = "UrlPre"}

How to use lists of keywords/phrases:

- Create a *.def file listing all gazetteers (with their major/minor categories):

```
cities.lst:location:city
```

```
organizations.lst:organization
```

```
surnames.lst:name:surname
```

```
forenames.lst:name:forename
```

- List one key phrase, name or keyword per line in each of these files, e.g. in cities.lst:

```
Abu Dhabi
```

```
Berlin
```

```
Chicago
```

```
Frankfurt
```

```
...
```

- Create a gazetteer Processing Resource in your application pipeline that references your *.def file.

Common Machine Learning Methods for Information Extraction

- **Hidden Markov Models (HMMs)** (recall from Week 2)
- Conditional Random Fields (CRF)
- Support Vector Machine (SVM)
- Artificial Neural Networks (NNs) (will be covered in Week 7)

BIO Encoding for Labels Classifying Text Spans (CoNLL)

- Annotate word tokens with inside/outside of class information:

The_0 Oracle_I-ORG CEO_0 's_0 name_0
is_0 Larry_I-PER Ellison_I-PER ._0
(horizontal)

or:

The 0
Oracle I-ORG
CEO 0
...
(vertical)

- B tag used to demarcate adjacent I tags (otherwise, two adjacent "I" tokens part of the same text span could not be distinguished from the case of two adjacent, separate text spans)

Annotation of Gold Data with BRAT

(<http://nlp.stanford.edu>)

Stanford CoreNLP

Output format: Visualise

Please enter your text here:

Chase Manhattan and its merger partner J.P.Morgan and Citibank, which was involved in moving about \$100 million for Raul Salinas de Gortari, brother of a former Mexican president, to banks in Switzerland, are also expected to sign on.

Submit Clear

Part-of-Speech:

Chase Manhattan and its merger partner J.P.Morgan and Citibank, which was involved in moving about \$100 million for Raul Salinas de Gortari, brother of a former Mexican president, to banks in Switzerland, are also expected to sign on.

Named Entity Recognition:

Chase Manhattan and its merger partner J.P.Morgan and Citibank, which was involved in moving about \$100 million for Raul Salinas de Gortari, brother of a former Mexican president, to banks in Switzerland, are also expected to sign on.

Coreference:

Chase Manhattan and its merger partner J.P.Morgan and Citibank, which was involved in moving about \$100 million for Raul Salinas de Gortari, brother of a former Mexican president, to banks in Switzerland, are also expected to sign on.

Basic dependencies:

Chase Manhattan and its merger partner J.P.Morgan and Citibank, which was involved in moving about \$100 million for Raul Salinas de Gortari, brother of a former Mexican president, to banks in Switzerland, are also expected to sign on.

Feature Extraction – Spam Classification Example (1/2)

- **Feature:** a piece of evidence intended to help the classifier map the input to the right target class
- **Feature vector:** a vector \vec{F} , the components $F_j = \phi_j(d_i)$, of which are results applying a feature function to the data point d_i
- Example: “Spam versus Ham” email?
 - number of “!”s included in email body
 - length of the email in characters
 - does the word “cash” occur in the title or body?
- Example feature vectors:
 - (2, 2392, no) \mapsto HAM (genuine e-mail)
 - (4, 520, yes) \mapsto SPAM
 - (1, 2392, no) \mapsto HAM
 - (0, 16337, no) \mapsto HAM
 - (1, 6i320, yes) \mapsto SPAM

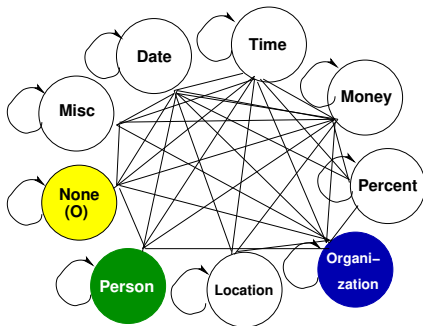
In English:

- Token begins with capital letter: `[A-Z] [a-z]+`
- Token ends in characters `-man`
- Token to the left is from list Mr., Mrs., Miss, Dr., Prof., Sir, Lord, CEO, ...
- Token to the right is academic title, affiliation: BSc, Ph.D., M.A., FRS, M.P.
- Tokens to the right are `[,]? who`
- Token is from a list (gazetteer) of names

- Determine the model's parameters from data: induction
- Training corpus has counts, from which model probabilities are computed
- Smoothing: re-distribution of probability mass from seen to unseen events (to avoid zero probabilities, which make any product go zero)

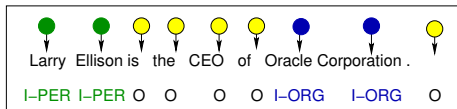
HMMs for IE e.g. in Nymble (Bikel et al., 1997)

HMM Transitions:



Not shown: initial tranistions/probabilities

Most likely state (label) sequence (Viterbi alg.):



Features in Nymble (Bikel et al., 1997)

Word Feature	Example Text	Intuition
twoDigitNum	90	Two-digit year
fourDigitNum	1990	Four digit year
containsDigitAndAlpha	A8956-67	Product code
containsDigitAndDash	09-96	Date
containsDigitAndSlash	11/9/89	Date
containsDigitAndComma	23,000.00	Monetary amount
containsDigitAndPeriod	1.00	Monetary amount, percentage
otherNum	456789	Other number
allCaps	BBN	Organization
capPeriod	M.	Person name initial
firstWord	<i>first word of sentence</i>	No useful capitalization information
initCap	Sally	Capitalized word
lowerCase	can	Uncapitalized word
other	,	Punctuation marks, all other words

Quantifying System Quality

- Basis of system quality is a quantitative comparative evaluation
- When comparing predicted class (system) with actual class, there are 4 cases:
 - TP (true positive)
 - FP (false positive)
 - FN (false negative)
 - TN (true negative)

- **Confusion matrix:**

TP	FP
FN	TN

Accuracy, Precision and Recall

- **Sample size:** Number of data points N (e.g. tokens) to be classified
- **Accuracy** (A, Correctness): $A = \frac{TP+TN}{TP+TN+FP+FN} = \frac{C}{N}$
is *not* a good measure to use
→ artificially high values (if label bias)
- **Precision** (P, Positive Predictive Value): $P = \frac{TP}{TP+FP}$
“fraction of retrieved instances that are relevant”
retrieved and correct over total retrieved (which proportion of responses were actually correct?)
- **Recall** (R, Sensitivity): $R = \frac{TP}{TP+FN}$
“fraction of relevant instances that are retrieved”
retrieved and correct over all correctly retrieved and missing
(which proportion of correct responses was actually retrieved?)

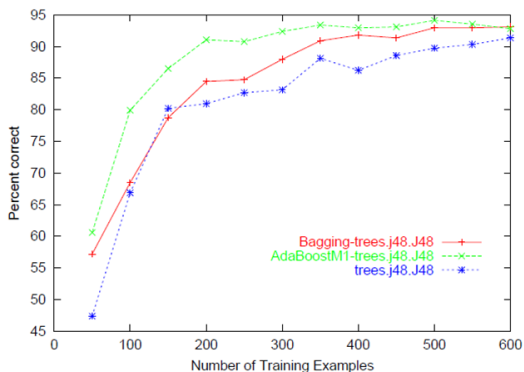
- Combined metric: **harmonic mean** of Precision and Recall
- Introduced by van Rijsbergen (1979)
- $F = (1 + \beta^2) \times \frac{PR}{\beta^2 P + R}$
- $F1 = F_{\beta=1} = \frac{2PR}{P+R}$ (equal weight of P and R)

Confusion Matrix (Source: Wikipedia)

- Most useful tool for performance evaluation

		predicted condition			
		total population	prediction positive	prediction negative	
true condition	condition positive	True Positive (TP)	False Negative (FN) (type II error)	True Positive Rate (TPR), Sensitivity, Recall, Probability of Detection $= \frac{\Sigma \text{ TP}}{\Sigma \text{ condition positive}}$	False Negative Rate (FNR), Miss Rate $= \frac{\Sigma \text{ FN}}{\Sigma \text{ condition positive}}$
	condition negative	False Positive (FP) (Type I error)	True Negative (TN)	False Positive Rate (FPR), Fall-out, Probability of False Alarm $= \frac{\Sigma \text{ FP}}{\Sigma \text{ condition negative}}$	True Negative Rate (TNR), Specificity (SPC) $= \frac{\Sigma \text{ TN}}{\Sigma \text{ condition negative}}$
		Accuracy $= \frac{\Sigma \text{ TP} + \Sigma \text{ TN}}{\Sigma \text{ total population}}$	Positive Predictive Value (PPV), Precision $= \frac{\Sigma \text{ TP}}{\Sigma \text{ prediction positive}}$	False Omission Rate (FOR) $= \frac{\Sigma \text{ FN}}{\Sigma \text{ prediction negative}}$	Positive Likelihood Ratio (LR+) $= \frac{\text{TPR}}{\text{FPR}}$
			False Discovery Rate (FDR) $= \frac{\Sigma \text{ FP}}{\Sigma \text{ prediction positive}}$	Negative Predictive Value (NPV) $= \frac{\Sigma \text{ TN}}{\Sigma \text{ prediction negative}}$	Negative Likelihood Ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$
					Diagnostic Odds Ratio (DOR) $= \frac{\text{LR}+}{\text{LR}-}$

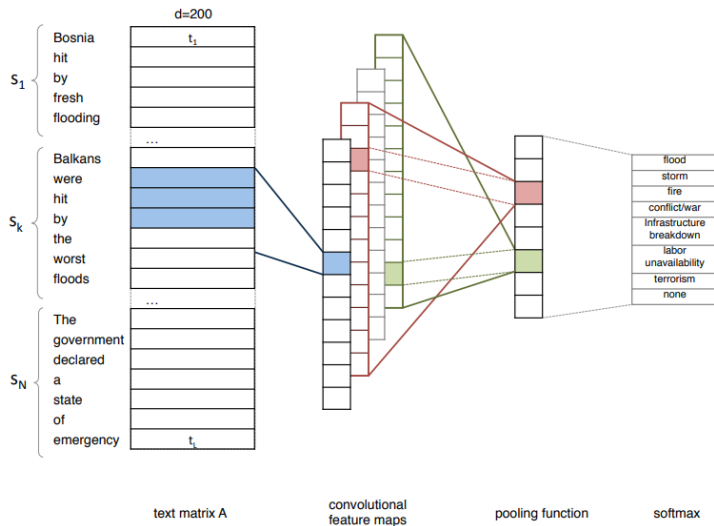
Learning Curve: Do We Have Enough Gold Data?



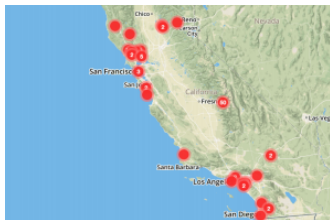
There are many application domains for IE systems:

- Bio-medical IE applications (genes & proteins)
- Financial IE applications (mergers & acquisitions, CEO changes))
- Legal IE applications (judges & attorneys)
- Intelligence & Police IE applications (terrorists & crimes)
- e-Commerce IE applications (brands & products)

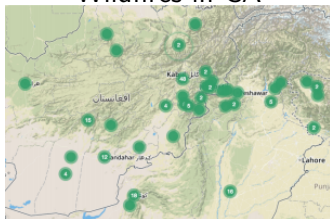
Event Extraction from News with Neural Models (Nugent et al., 2017) (1/2)



Event Extraction from News with Neural Models (Nugent et al., 2017) (2/2)



Wildfires in CA



Conflicts in Afghanistan

Current Developments & Future Directions in IE

- Utilizing robust syntactic and semantic components where available
- From classic IE to *open IE*
- From small data sets to Web-scale data sets (bigger data and a simpler algorithm usually beats smaller data and a more sophisticated algorithm!)
- Knowledge-rich methods are going to be combined with machine learning methods
- Use of word embeddings and deep (= multiple hidden layers) neural networks,

In this session, we learned:

- what information extraction is, and a bit of its history
- the difference between rule-based and machine learning approaches
- how named entities can be extracted from text

References (1/2)

- Bird/Klein/Loper (2009), *Natural Language Processing with Python*, Sebastopol, CA: O'Reilly
- Costantino and Coletti (2008), *Information Extraction in Finance*, Southampton: WIT Press
- Cunningham et al. (2013), *Developing Language Processing Components with GATE Version 7 (a User Guide)*, University of Sheffield, [online] <http://gate.ac.uk/sale/tao/>, Chapter 2-3, 5-6, and 10.
- Hastie, Tibshirani and Friedman (2009), *The Elements of Statistical Learning (2nd ed.)*, New York, NY: Springer
- Jurafsky/Martin (2008), *Speech and Language Processing (2nd ed.)*, Upper Saddle River, NJ: Prentice Hall
- Manning/Schütze (1999) *Statistical Natural Language Processing* MIT Press.
- Moens (2006), *Information Extraction: Algorithms and Prospects in a Retrieval Context* Dordrecht: Springer

- Nugent & Leidner (2016), “Risk Mining: Company-Risk Identification from Unstructured Sources” *Proc. ICDM*
- Plachouras, Leidner & Garrow (2016) “Quantifying Self-Reported Adverse Drug Events on Twitter: Signal and Topic Analysis” *Ann. Meeting of the Social Media Soc.*
- Pustejovsky/Stubbs (2012), Natural Language Annotation for Machine Learning, Sebastopol, CA: O'Reilly, Chapters 5-6 and 8
- Zanasi (2005), *Text Mining and its Applications to Intelligence, CRM and Knowledge Management*, Southampton: WITPress

Backup Slides

Some Well-Known IE Systems

SAM	PAM	FRUMP	ATRANS	Proteus
PLUM	Tabula Rasa	LOLITA	Scrabble	NameTag
Alembic	Hasten	NLToolset	IE ²	TurboTag
SIFT	TASC	FACILE	AutoSlog	ANNIE
LasIE	Elie	NetOwl	Oki	BALIE
Diderot	JASPER	Tacitus	FASTUS	SCISOR
Circus	REES	Identifinder	Nymble	InfoXtract

Homework: Constructing IE Pipelines with GATE (1/2)

- File→New Corpus Pipeline, give name e.g. TestIE
- Processing Resources→Add Annie
- Processing Resources→Add Document Reset, ANNIE Sentence Splitter, ANNIE English Tokeniser, ANNIE Gazetteer, and ANNIE NE Transducer.
- Double click on TestIE and move the components Document Reset, ANNIE Sentence Splitter, ANNIE English Tokeniser, ANNIE Gazetteer, and ANNIE NE Transducer from the left list to the right using ">>"; ensure the order is exactly as given here from top to bottom!
- Language Resources→Add Corpus Document, click on sourceURL and select a text file to process

Homework: Constructing IE Pipelines with GATE (2/2)

- Language Resources→Add Corpus, then add document created above via drop-down menu to the empty corpus
- Double-click on TestIE Application, select corpus in drop-down and click Run this Application
- Double-click on your test document under Language Resources and select the annotations (PERSON, ORGANIZATION) you would like to view

Installing NLTK

```
$ sudo apt-get install python-nltk
$ python
>>> import nltk
>>> nltk.download()
```

Using NLTK – A KWIC Concordance

Let's extract a Keyword-in-Context concordance of the word "grail" in Monty Python and the Holy Grail:

```
from nltk.book import *  
text6.concordance("grail")
```

```
ARTHUR : If you will not show us the grail , we shall take your  
s required if the quest for the Holy grail were to be brought to  
should separate , and search for the grail individually . [ clop  
AD : You are the keepers of the Holy grail ? ZOOT : The what ? G  
il ? ZOOT : The what ? GALAHAD : The grail . It is here . ZOOT :  
ease ! In God ' s name , show me the grail ! ZOOT : Oh , you hav  
rment me no longer . I have seen the grail ! PIGLET : There ' s  
en the Grail ! PIGLET : There ' s no grail here . GALAHAD : I ha  
are you going ? GALAHAD : I seek the grail ! I have seen it , he  
which , I have just remembered , is grail - shaped . It ' s not  
blem . GALAHAD : It ' s not the real grail ? DINGO : Oh , wicked
```


Using NLTK – A POS Tagging Example

```
>>> import nltk
>>> sentence = """Although big data is a cool topic, Monday
morning lectures are never easy."""
>>> tokens = nltk.word_tokenize(sentence)
>>> tokens

>>> tagged = nltk.pos_tag(tokens)
['Although', 'big', 'data', 'is', 'a', 'cool', 'topic', ',',
 'Monday', 'morning', 'lectures', 'are', 'never', 'easy',
 '.']

>>> tagged
[('Although', 'IN'), ('big', 'JJ'), ('data', 'NNS'), ('is',
'VBZ'), ('a', 'DT'), ('cool', 'NN'), ('topic', 'NN'), (',',
','), ('Monday', 'NNP'), ('morning', 'NN'), ('lectures',
'NNS'), ('are', 'VBP'), ('never', 'RB'), ('easy', 'JJ'),
('.', '.')]

```

Using NLTK – the Built-in Named Entity Tagger

```
>>> nltk.ne_chunk(nltk.pos_tag(nltk.word_tokenize(
    "President Obama had to face the decision
    whether to visit Putin or not.")))
Tree('S', [('President', 'NNP'), Tree('PERSON', [('Obama',
'NNP')]), ('had', 'VBD'), ('to', 'TO'), ('face', 'VB'),
('the', 'DT'), ('decision', 'NN'), ('whether', 'IN'),
('to', 'TO'), ('visit', 'VB'), Tree('PERSON', [('Putin',
'NNP')]), ('or', 'CC'), ('not', 'RB'), ('.', '.')])])
```

Using NLTK – A Regular Expression Tagger

```
>>> regexp_tagger = nltk.RegexpTagger(  
...     [(r'^-?[0-9]+(.[0-9]+)?$', 'CD'),# cardinal numbers  
...     (r'(The|the|A|a|An|an)$', 'AT'),# articles  
...     (r'.*able$', 'JJ'),             # adjectives  
...     (r'.*ness$', 'NN'),             # nouns formed from adjectives  
...     (r'.*ly$', 'RB'),               # adverbs  
...     (r'.*s$', 'NNS'),               # plural nouns  
...     (r'.*ing$', 'VBG'),             # gerunds  
...     (r'.*ed$', 'VBD'),             # past tense verbs  
...     (r'.*', 'NN')                   # nouns (default)  
... ])  
>>> regexp_tagger.tag(test_sent)  
[('The', 'AT'), ('Fulton', 'NN'), ('County', 'NN'), ('Grand',  
  'NN'), ('Jury', 'NN'), ('said', 'NN'), ('Friday', 'NN'),  
  ('an', 'AT'), ('investigation', 'NN'), ('of', 'NN'),  
  ("Atlanta's", 'NNS'), ('recent', 'NN'), ('primary', 'NN'),  
  ('election', 'NN'), ('produced', 'VBD'), (''', 'NN'),  
  ('no', 'NN'), ('evidence', 'NN'), ....]
```

Using NLTK – Drawing Analysis Trees

```
from nltk.tree import *
from nltk.draw import tree
entities = nltk.ne_chunk(nltk.pos_tag(
    nltk.word_tokenize("Maria likes fast computers.")))
entities.draw()
```



Sequence Tagging

- It is one thing to classify a data point in isolation

Example (1): “dinosaur” \mapsto noun

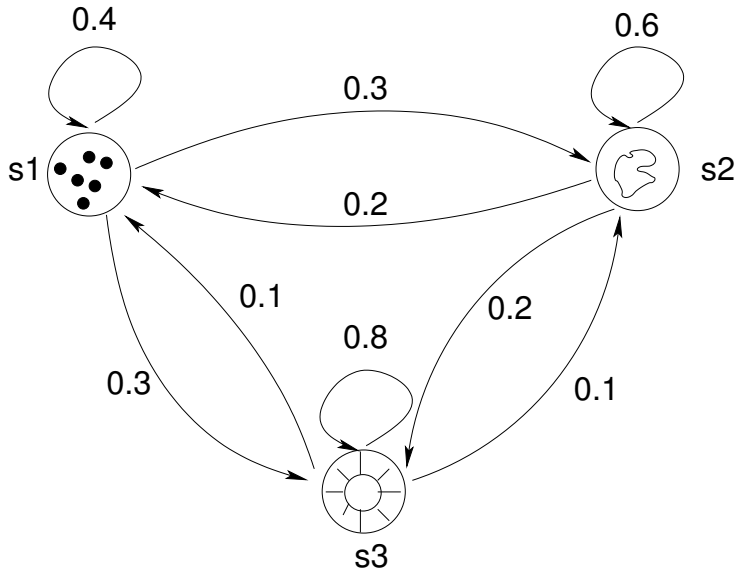
- It is quite another thing to classify, taking into account *context*

Example (2a): “(let’s) walk (there)” \mapsto walk/VVB

Example (2b): “(a) walk (in the rain)” \mapsto walk/NN

- If “walk” has two possible readings, how can we compute the right one? \rightarrow **Ambiguity**
- Example (3): I can can the can . \mapsto
I/PRP can/AUX can/VVB the/AT can/NN ./.

Markov Model – Introductory Example (1/2)



Markov Model – Introductory Example (2/2)

- Simple probabilistic finite-state model of the weather (**Markov Chain**)
- 3 States: s_1 : rainy, s_2 : cloudy, s_3 : sunny
- Transitions labelled with probabilities $a_{ij} \geq 0, \forall i, j \leq N$,
 $\sum_{j=1}^N a_{ij} = 1, \forall i$
- What is the probability of the observation sequence $O =$
(*sunny, sunny, sunny, rainy, rainy, sunny, cloudy, sunny*)?

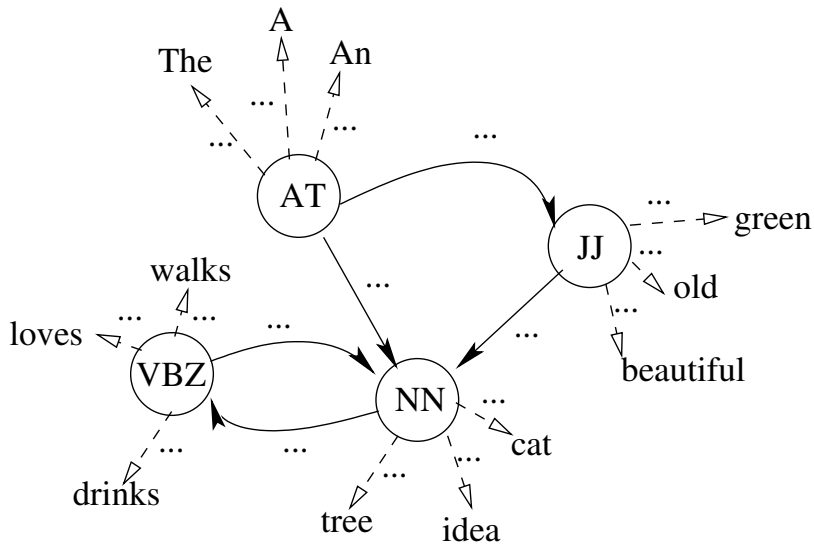
$$\begin{aligned}P(O|Model) &= P[s_3, s_3, s_3, s_1, s_1, s_3, s_2, s_3|Model] \\&= P[3]P[3|3]^2P[1|3]P[1|1]P[3|1]P[2|3]P[3|2] \\&= \pi_3 a_{33}^2 a_{31} q_{11} a_{11} a_{13} a_{32} a_{23} \\&= 1.0 \cdot 0.8^2 \cdot 0.1 \cdot 0.4 \cdot 0.3 \cdot 0.1 \cdot 0.2 \\&= 1.536 \times 10^{-4}\end{aligned}$$

Hidden Markov Models (HMM)

A **Hidden Markov Model** $\lambda = (Q, \Sigma, A, B, \Pi)$ is formally defined as

- a set of N *hidden states* $Q = (q_1, \dots, q_N)$, $N = |Q|$
- a set Σ of M *observation symbols*, $M = |\Sigma|$
- a state-transition distribution (*transition probabilities*)
 $A = \{a_{ij} = P(q_{t+1} = j | q_t = i)\}, 1 \leq i, j \leq N$
- an observation symbol probability distribution (*emission probabilities*)
 $B = \{b_{ik} = b_o(o_k) = P(o_k | q_i)\}, 1 \leq k \leq M, o_k \in \Sigma$ (the probability that the output is o_k , given that the current state is q_i)
- an *initial state distribution*
 $\Pi = \{\pi_i = P(q_i | t = 1)\}, 1 \leq i \leq N$

Hidden Markov Models (HMM) – POS Tagging Example



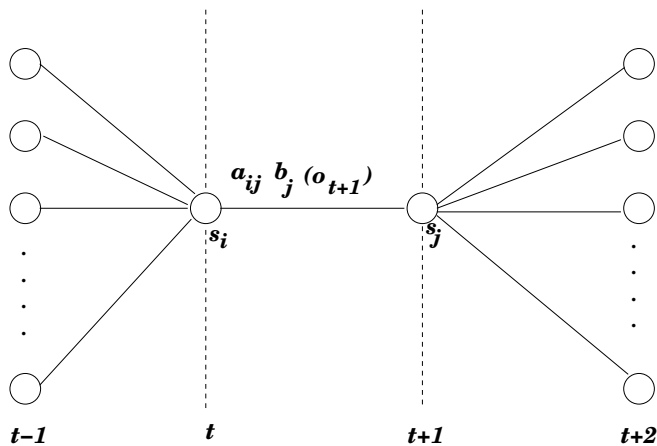
Three Questions to an HMM

- **Problem 1 (evaluation).** Given the observation sequence $O = (o_1, o_2, \dots, o_T)$ and a model $\lambda = (Q, \Sigma, A, B, \Pi)$, how do we efficiently compute $P(O|\lambda)$, the probability of a observation sequence, given a model?
- **Problem 2 (decoding).** Given the observation sequence $O = (o_1, o_2, \dots, o_T)$ and the model λ , how can we find a corresponding state sequence $q^* = q_1, q_2, \dots, q_T$ that most likely generated O (“optimally explains the observations”)?
- **Problem 3 (learning).** How do we estimate the model parameters A , B and Π so as to maximize $P(O|\lambda)$?

The Viterbi algorithm (1/2) (Viterbi, 1967)

- Finding the most likely hidden state sequence, given an observation, requires at first glance exponential runtime complexity ($O(N^n)$), as all possible states are valid hypotheses for each observation, so candidate states/readings multiply)
- Luckily, we can find a linear-runtime ($O(n)$) solution using dynamic programming (remembering partial solutions at each step) in a data structure called **trellis**.
- The Viterbi algorithm uses one pass from left to right looking at initial, transition and emission probabilities while storing a history of the locally “best” (most likely) state (back pointer).

Trellis Data Structure



The Viterbi algorithm (2/2) (Viterbi, 1967)

Let $\delta_t(i) = \max_{q_1, \dots, q_{t-1}} P[q_1 q_2 \dots q_{t-1}, q_t = i, o_1 o_2 \dots o_t | \lambda]$.

- 1 Initialization. (Start with initial state probabilities)

$$\delta_1(i) = \pi_i b_i(o_1), 1 \leq i \leq N$$

$$\phi_1(i) = 0$$

- 2 Recursion. (main part, see previous slide)

$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij} b_j(o_t)], 2 \leq t \leq T, 1 \leq j \leq N$$

$$\phi_t(j) = \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}], 2 \leq t \leq T, 1 \leq j \leq N$$

- 3 Termination.

$$P^* = \max_{1 \leq i \leq N} [\delta_T(i)]$$

$$q_T^* = \arg \max_{1 \leq i \leq N} [\delta_T(i)]$$

- 4 Path backtracking (finding the most likely hidden state sequence).

$$q_t^* = \phi_{t+1}(q_{t+1}^*), t = T-1, T-2, \dots, 1$$

JRC European Media Monitor (EMM) (1/3) – Steinberger et al.



JRC European Media Monitor (EMM) (2/3)



EMM NewsBrief

EMM NewsExplorer

News Explorer

RSS feed for this entity

Daily News Analysis, across languages and over time

News Summary

About EMM NewsExplorer

News language and date

Language or country:
en - English

Date:
Nov 2013

Mo	Tu	We	Th	Fr	Sa	Su
				1	2	3
4	5	6	7	8	9	10
11	12	13	14	15	16	17
18	19	20	21	22	23	24
25	26	27	28	29	30	

Analysis over time

Timeline

Timeline [en] for 11/2013



Osama bin Laden

Information about this person was last updated on Saturday, November 16, 2013.

Names

Osama bin Laden (Eu,tr)
Bin Laden (da,tr)
bin Laden (da,tr)
Osama Bin Laden (da,sw)
بن لادن (ar,fa)
Oussama ben Laden (fr)
Ben Laden (de,ro)
Usama bin Ladin (de,tr)
Oussama Ben Laden (tr)
بن لادن (ar)
اسامة بن لادن (ar)
Osama bin Laden (bg,sr)
Osama ben Laden (es,ro)
Osama bin Ladhna (sl)
Bin Ladin (de,tr)
Usame Bin Ladin (tr)
Бин Ладен (bg,ru)
Usama bin Laden (de,sv)
Ben Laden (bg,ru)
Osama bin-Laden (da,pt)
Usame bin Ladin (tr)
بن لادن (ar)
ben Laden (fr,ro)

Key Titles and Phrases

al Qaeda leader (en - 284)
al-Qaeda leader (en - 204)
lider de Al Qaeda (es,pt - 126)
Al Qaeda leader (en - 94)
chef (da,fr - 546)
terrorchef (de - 101)
Al-Qaeda leader (en - 81)
terroristenführer (de - 75)
lider da Al Qaeda (pt - 62)
leader (en,it - 417)
terroristenleider (nl - 60)
führer (de - 103)
leider (nl - 157)
Al-Qaeda chief (en - 39)
al Qaeda chief (en - 36)
al-Qaeda chief (en - 30)
terroristenführers (de - 35)
lider da Al-Qaeda (pt - 32)
Al-Qaida leader (en - 30)
anführer (de - 63)
lider (es,pt - 103)
terroristenchef (de - 27)
born (en - 72)

External resources



Image obtained automatically from Wikipedia

Photo Wikipedia entry

Explore Relations



Related People

Barack Obama (447)
John Kerry (198)
Ayman al-Zawahiri (138)
George W. Bush (122)
Hillary Rodham Clinton (105)
Bashar Assad (101)
Kathryn Bigelow (78)
Hamid Karzai (75)
Chuck Hagel (75)
John McCain (74)
Nawaz Sharif (70)
Bill Clinton (64)
François Hollande (61)
Jay Carney (60)
David Cameron (60)
Vladimir Putin (59)
Steven Spielberg (58)
Ben Affleck (57)
Bradley Manning (56)
Edward Snowden (56)
Mohammed Morsi (55)
Daniel Day-Lewis (54)
Joseph Biden (53)
Leon Panetta (52)
Julian Assange (50)
Susan Rice (50)

Latest Clusters - English

[fr] [ar] [bg] [ru] [pt] [de] [hu] [da] [nl] [es] [no] [sv] [it] [tr] [ro] [sw] [et] [sl] [fa]

On the rise: Polio cases pass 2012 total
tribune 13-NOV-13

Blasts at Indian political rally kill 5
cnn 27-OCT-13

PM optimistic about investments in Pakistan
thehindu 13-NOV-13

Afghan president criticises timing of death of Pakistan
Taliban leader
telegraph 04-NOV-13

Sartaj Aziz says US understands Pakistan viewpoint on
drone attacks
TheFrontierPost 25-OCT-13

President Obama welcomes PM Sharif at White House as
thehindu 13-NOV-13

JRC European Media Monitor (EMM) (3/3)



Zoom:



Pan:



Search for a name:

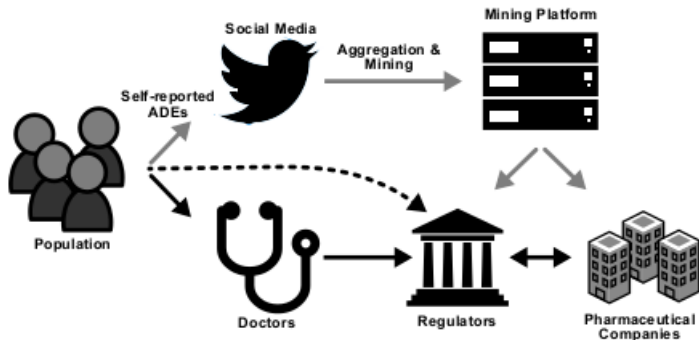
S

Search results:

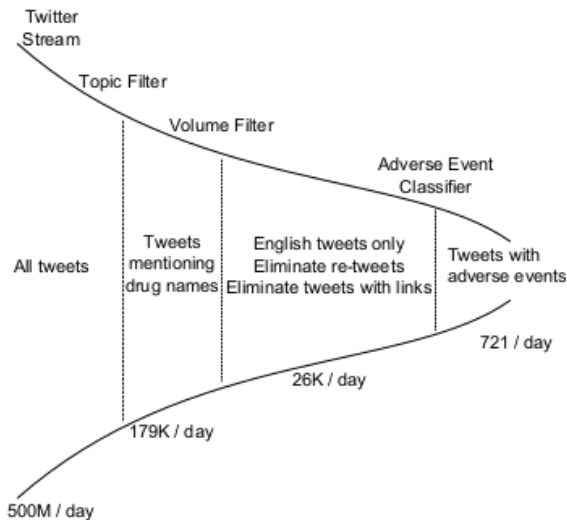
Click on a name to add it to the relation network:

These persons or organisations are already shown in the graphic:

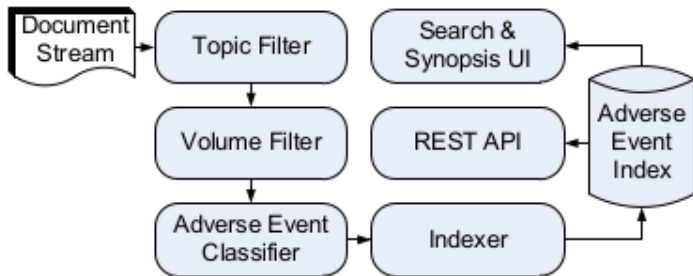
Pharmacology IE – Medical Drug Use Monitoring (Plachouras et al., 2016) (1/6)



Pharmacology IE – Medical Drug Use Monitoring (Plachouras et al., 2016) (2/6)



Pharmacology IE – Medical Drug Use Monitoring (Plachouras et al., 2016) (3/6)



Pharmacology IE – Medical Drug Use Monitoring (Plachouras et al., 2016) (4/6)

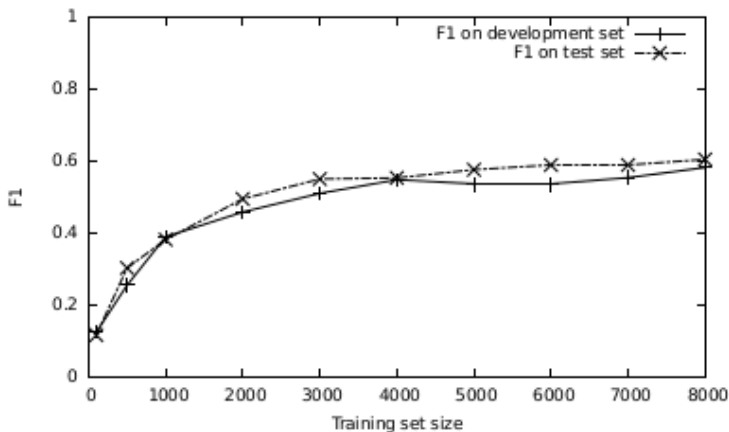


Pharmacology IE – Medical Drug Use Monitoring (Plachouras et al., 2016) (5/6)

	Tweet text
TP	wellbutrin doesnt even work for me it just makes me really anxious idk why im still taking it
TP	Took some ibuprofen that has me so drowsy
TP	Insomnia and heart palpitations due to prednisone. Takteng mga side effects 'to. /wrist
TP	This Vicodin is making me feel like a tweaker because I'm so itchy!
TP	guys i took an ibuprofen 2 days ago and i still have heart palpitations
TN	I once gave a treatment to some patients with glutathione which I knew later clinic use it. Most of them have nausea for 5 mins after inj.
FN	I seriously never have any energy thanks accutane lol @probs_accutane all I want to do is sleep
FP	My mouth taste like 1200 Mg's of ibuprofen yet my head still hurts and I'm still feeling dizzy. Wtf
FP	@ash_hein they gave me Tylenol 3s & yea kinda my mouth still hurts a little & I'm still swollen

Pharmacology IE – Medical Drug Use Monitoring (Plachouras et al., 2016) (6/6)

Model	C	w^+	Development			Test			
			P	R	F1	P	R	F1	
Baseline			0.269	0.683	0.386	0.267	0.705	0.387	
BIN_NGRAM1,2	0.050	8	0.636	0.504	0.562	0.560	0.540	0.549	*
ALL FEATURES	0.025	9	0.573	0.590	0.582	0.550	0.669	0.604	*†



Where to Get Linguistic Gold Data From

- Search the Linguistic Data Consortium (LDC) catalog (<http://catalog.ldc.upenn.edu>)
- Search ELRA catalog (<http://catalog.elra.info>, <http://www.hlt-evaluation.org>)
- Browse DFKI's LT World <http://www.lt-world.org>
- Browse MetaNet (<http://www.meta-net.eu>)
- Browse the ACL Anthology (<http://acl.ldc.upenn.edu>), ACM Digital Library (<http://dl.acm.org>) and IEEEExplore (<http://ieeexplore.ieee.org>)
 - some papers on a topic may contain associated data
 - some conferences, notably the bi-annual LREC, specialize on presenting new linguistic resources
- Contact expert researchers
- Recruit linguists/domain experts
- Crowdsourcing